



# RELIABLE EXPLAINABLE AI FRAMEWORK FOR SOFTWARE DEFECT PREDICTION

Submitted in partial fulfilment of the requirements of the degree of  
**Master of Technology**

by

**Udit Jain**  
(Roll No: 24CSM1R23)

Supervisor:

**Dr. Manjubala Bisi**  
*Assistant Professor, Department of CSE*

Department of Computer Science and Engineering  
National Institute of Technology, Warangal  
India

**February 05, 2026**

# Acknowledgement

I would like to express my sincere gratitude to my project supervisor,

**Dr. Manjubala Bisi**, Assistant Professor, Department of Computer Science and Engineering, National Institute of Technology Warangal, for her valuable guidance, encouragement, and continuous support throughout the course of this project. Her insights, suggestions, and constructive feedback have been instrumental in shaping this work.

I also extend my heartfelt thanks to all the faculty members and staff of the Department of Computer Science and Engineering, NIT Warangal, for providing an excellent academic and research environment, and for facilitating the resources required for the successful completion of this project.

I am deeply grateful to my friends and batchmates for their ideas, encouragement, and assistance during various phases of this work. Their moral support played a key role in keeping me motivated.

A special note of appreciation goes to my family for their unconditional support, patience, and constant encouragement throughout my academic journey.

Finally, I am thankful to the National Institute of Technology Warangal for providing the infrastructure and learning environment that enabled me to carry out this project successfully.

**Udit Jain**

Roll No: 24CSM1R23

# Declaration

I hereby declare that the work presented in this dissertation is the outcome of my own research carried out under the supervision of

**Dr. Manjubala Bisi**, Assistant Professor, Department of Computer Science and Engineering, National Institute of Technology, Warangal.

To the best of my knowledge, this dissertation has not been submitted either in part or in full for the award of any other degree or diploma at this or any other institution. Wherever the work of others has been used, it has been duly acknowledged and properly cited.

I further affirm that I have maintained the principles of academic integrity and honesty throughout the preparation of this dissertation, and I take full responsibility for the content presented herein.

(Signature)

**Udit Jain**

Roll No: 24CSM1R23

# Approval Sheet

This is to certify that the dissertation work entitled

**”Reliable Explainable AI Framework for Software Defect Prediction”**

submitted by **Udit Jain (Roll No: 24CSM1R23)**

has been examined and is hereby approved for the award of the degree of  
**Master of Technology.**

**Examiners:**

Prof. Manish Kumar Bajpai

Prof. Earnest Paul Ijjina

Prof. Sarath Babu

Prof. Sanjaya Kumar Panda

**Supervisor:**

Dr. Manjubala Bisi

**Head of Department (CSE):**

Dr. Rashmi Ranjan Rout

**Department of Computer Science and Engineering**

National Institute of Technology, Warangal, India

# Certificate

This is to certify that the dissertation work entitled  
**”Reliable Explainable AI Framework for Software Defect Prediction”**

is a bonafide record of research carried out by **Udit Jain (Roll No: 24CSM1R23)** under my supervision.

The dissertation has been submitted to the Department of Computer Science and Engineering, National Institute of Technology, Warangal, in partial fulfilment of the requirements for the award of the degree of **Master of Technology** during the academic year 2024–2025.

(Signature)

**Dr. Manjubala Bisi**

Assistant Professor

Department of Computer Science and Engineering

National Institute of Technology, Warangal

# Abstract

Software defect prediction plays a crucial role in improving software quality by identifying defect-prone modules early in the development cycle. While machine learning models such as Random Forests and other ensemble techniques have shown promising results in predicting software defects, their adoption in practice is often hindered by a lack of trust in model explanations. Existing studies primarily focus on accuracy and interpretability but rarely assess the reliability of these explanations.

This dissertation proposes a **Reliable Explainable AI Framework** that integrates SHAP (SHapley Additive exPlanations) with novel reliability metrics. In addition to standard measures of generalizability, concordance, and stability, the framework introduces three behavior-grounded metrics: (i) *Gradient-of-Loss Reliability (GLR)* to evaluate alignment between SHAP feature importance and true loss sensitivity, (ii)  $\epsilon$ -*Perturbation Hit@k* to test whether SHAP-identified top features genuinely affect model loss under small perturbations, and (iii) *Expected Calibration Error (ECE)* to assess the reliability of predicted probabilities.

Experiments were conducted on the NASA PROMISE CM1 dataset using Random Forest classifiers. Results demonstrate that while conventional SHAP explanations achieve high generalizability, their reliability varies across folds and conditions. The proposed framework provides a composite reliability score, offering a more comprehensive and trustworthy evaluation of explainability.

This work bridges the gap between interpretability and reliability in explainable AI for software defect prediction. By validating explanations against model behavior and calibration, the framework enhances confidence for developers and stakeholders, moving closer to trustworthy AI in software engineering.

**Keywords:** Explainable AI, Reliability, SHAP, Software Defect Prediction, Random Forest, GLR, Hit@k, Calibration

# Contents

<b>Acknowledgement</b>	<b>i</b>
<b>Declaration</b>	<b>ii</b>
<b>Certificate</b>	<b>iv</b>
<b>Abstract</b>	<b>v</b>
<b>List of Figures</b>	<b>viii</b>
<b>List of Abbreviations</b>	<b>x</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Context and Motivation .....	1
1.2 Challenges in Current Systems.....	2
1.3 Reliability as a Core Requirement.....	2
1.4 Objectives and Scope of the Dissertation.....	3
1.5 Expected Contributions.....	3
1.6 Organization of the Report .....	3
<b>2 Background</b>	<b>5</b>
2.1 Reliability in Explainable AI.....	5
2.2 Behavior-Grounded Reliability .....	6
2.3 Machine Learning for Defect Prediction.....	6
2.4 Explainability as a Tool for Reliability .....	6
2.5 Summary .....	7
<b>3 Related Work</b>	<b>8</b>
3.1 Machine Learning for Software Defect Prediction .....	8
3.2 Explainable AI in Software Engineering.....	8
3.3 Reliability of Explanations .....	9
3.4 Positioning in Recent Literature (2023–2025) .....	10
3.5 Gaps and Open Questions.....	10
3.6 Need for This Study .....	10
3.7 Summary .....	11
<b>4 Proposed Solution</b>	<b>12</b>
4.1 System Overview .....	12

4.2	Dataset and Preprocessing .....	14
4.3	Model Training and Standard Evaluation .....	16
4.4	Explainability using SHAP .....	17
4.5	Reliability Assessment .....	17
4.5.1	Classical Reliability Metrics .....	17
4.5.2	Behavior-Grounded Reliability Metrics .....	18
4.5.3	Composite Reliability Score .....	18
4.6	Pipeline Pseudocode .....	19
4.7	Visualization Outputs .....	19
4.8	Summary .....	21
<b>5</b>	<b>Experimental Results</b> .....	<b>22</b>
5.1	Overall Experimental Setup .....	22
5.2	Same Project Evaluation (Train = Test) .....	22
5.2.1	Baseline Results (Same Project) .....	23
5.2.2	Proposed Results (Same Project) .....	23
5.2.3	WSR Results (Same Project) .....	23
5.2.4	Conclusion (Same Project) .....	23
5.3	Cross Project Evaluation (Train $\neq$ Test) .....	24
5.3.1	WSR Results (Cross Project) .....	24
5.3.2	Conclusion (Cross Project) .....	25
5.4	Summary .....	25
<b>6</b>	<b>Conclusion and Future Work</b> .....	<b>26</b>
6.1	Conclusion .....	26
6.2	Implications of Findings .....	27
6.3	Limitations .....	28
6.4	Future Work .....	28
6.5	Final Remarks .....	30
	<b>References</b> .....	<b>31</b>



# List of Figures

4.1	Proposed framework pipeline: dataset preparation, preprocessing, Random Forest training, SHAP-based explanations, and reliability evaluation.....	13
4.2	Top-10 influential software metrics for the cross-project pair <code>pc4.arff</code> $\rightarrow$ <code>pc1.arff</code> (AdaBoost), based on the <i>mean importance</i> computed on the test set. Larger bars indicate higher average contribution magnitude.....	19

# List of Tables

5.1	Baseline results in the Same Project setting (model-wise means with overall column means). .....	23
5.2	Proposed results in the Same Project setting (model-wise means with overall column means). .....	23
5.3	Wilcoxon Signed-Rank (WSR) test results for the Same Project setting (Baseline vs Proposed). .....	23
5.4	Fold-wise evaluation results under the Baseline approach (AdaBoost: Train = <code>pc4</code> , Test = <code>pc1</code> ). .....	24
5.5	Fold-wise evaluation results under the Proposed approach (AdaBoost: Train = <code>pc4.arff</code> , Test = <code>pc1.arff</code> ). .....	24
5.6	Wilcoxon Signed-Rank (WSR) test results comparing Baseline vs Proposed reliability scores on cross-project dataset pairs using AdaBoost. ....	25

# List of Abbreviations

<b>AI</b>	Artificial Intelligence
<b>ML</b>	Machine Learning
<b>DL</b>	Deep Learning
<b>SDP</b>	Software Defect Prediction
<b>XAI</b>	Explainable Artificial Intelligence
<b>SHAP</b>	SHapley Additive exPlanations (feature attribution method)
<b>LIME</b>	Local Interpretable Model-Agnostic Explanations
<b>RF</b>	Random Forest (tree-ensemble classifier)
<b>SVM</b>	Support Vector Machine
<b>CNN</b>	Convolutional Neural Network
<b>LCNN</b>	Lightweight Convolutional Neural Network
<b>MLP</b>	Multi-Layer Perceptron
<b>NAM</b>	Neural Additive Model (inherently interpretable)
<b>XNAM</b>	eXplainable Neural Additive Model (NAM with interactions)
<b>SMOTE</b>	Synthetic Minority Over-sampling Technique (imbalance handling)
<b>CV</b>	Cross-Validation
<b>AUC</b>	Area Under the ROC Curve (discrimination)
<b>F1</b>	Harmonic mean of Precision and Recall
<b>ECE</b>	Expected Calibration Error (probability calibration)
<b>GLR</b>	Gradient-of-Loss Reliability (behavior-grounded faithfulness)
<b>Hit@k</b>	Overlap of explainer top- $k$ features with top- $k$ loss-sensitive features
<b>Brier</b>	Brier Score (mean squared error of probabilities)
<b>LOC</b>	Lines of Code
<b>G/C/S</b>	Generalizability / Concordance / Stability (classical reliability)

# Chapter 1

## Introduction

Software systems form the backbone of modern society, supporting critical infrastructure, business operations, and everyday applications. As reliance on software grows, ensuring its quality and reliability has become increasingly important. Software defects not only raise maintenance costs and reduce performance but can also lead to significant financial losses or safety risks. Detecting defect-prone modules early in the development cycle enables developers to optimize resources, reduce testing overhead, and improve overall software quality.

### 1.1 Context and Motivation

Machine Learning (ML) models such as Random Forests, Gradient Boosting, and Neural Networks have demonstrated strong predictive power in software defect prediction (SDP). They can capture complex, non-linear relationships among software metrics such as size, complexity, or coupling. However, despite their predictive strength, these models are often treated as “black boxes,” making it difficult for practitioners to understand or trust their outputs.

Explainable AI (XAI) has emerged to address this issue by providing post-hoc explanations for model predictions. Among these techniques, SHAP (SHapley Additive exPlanations) is widely adopted for attributing feature importance in both local and global contexts. Yet, interpretability alone is not enough — what truly matters is whether these explanations are **reliable**. Explanations that vary across folds, disagree with model behavior, or fail to generalize to unseen data may create false confidence and hinder adoption in practice.

## 1.2 Challenges in Current Systems

Despite advancements in AI-driven defect prediction, several critical challenges remain:

- **Opaque Predictions:** High-performing models provide limited insight into the reasoning behind their outputs.
- **Data Constraints:** Software defect datasets are often imbalanced, noisy, and project-specific, affecting robustness.
- **Unreliable Explanations:** Post-hoc methods like SHAP or LIME are rarely evaluated for stability, generalizability, or alignment with the model’s actual behavior.
- **Trust Deficit:** Developers hesitate to rely on AI models if explanations cannot be validated or shown to be consistent under perturbations.

## 1.3 Reliability as a Core Requirement

Interpretability must be complemented by reliability for explanations to be meaningful. Reliable explanations should remain stable across data splits, align with internal model signals, and reflect genuine sensitivity of the model’s loss function to feature perturbations. To capture these aspects, reliability can be assessed through two complementary layers:

- **Classical Reliability Dimensions:** Generalizability (train–test consistency), Concordance (agreement with model-internal importances), and Stability (consistency across folds).
- **Behavior-Grounded Metrics:**
  - **Gradient-of-Loss Reliability (GLR):** Evaluates whether SHAP ranks align with loss sensitivity derived from perturbations.
  - **$\epsilon$ -Perturbation Hit@k:** Validates whether SHAP’s top- $k$  features also cause the largest changes in model loss under small perturbations.
  - **Expected Calibration Error (ECE):** Measures how well predicted probabilities reflect actual defect frequencies, ensuring probability honesty.

These metrics go beyond interpretability and provide quantitative checks to ensure that explanations are not only understandable but also trustworthy.

## 1.4 Objectives and Scope of the Dissertation

The objective of this dissertation is to design a **Reliable Explainable AI Framework** for software defect prediction. The scope of the work includes:

- Preprocessing defect datasets and addressing imbalance using methods such as SMOTE and class weighting.
- Training ML models, with Random Forests as the primary classifier, to distinguish defective from non-defective modules.
- Generating feature-level explanations using SHAP.
- Extending evaluation to reliability metrics including GLR,  $\epsilon$ -Hit@k, and ECE, alongside classical measures of generalizability, concordance, and stability.
- Proposing a composite reliability score that integrates these dimensions for a holistic assessment of explanation trustworthiness.

## 1.5 Expected Contributions

The contributions of this dissertation are expected to be:

- A structured framework that evaluates both interpretability and reliability of XAI in software defect prediction.
- Integration of behavior-grounded metrics (GLR, Hit@k, and ECE) into the reliability evaluation process.
- Empirical validation of the framework on NASA PROMISE datasets, demonstrating its effectiveness in producing stable and trustworthy explanations.
- Practical insights for software engineers and project managers to adopt AI-based defect prediction systems with greater confidence.

## 1.6 Organization of the Report

The rest of this dissertation is organized as follows:

- **Chapter 1: Introduction** — Describes motivation, challenges, objectives, and contributions.
- **Chapter 2: Background** — Reviews defect prediction, explainability, and reliability dimensions in XAI.

- **Chapter 3: Literature Survey** — Analyzes prior research, highlighting the gap in reliability assessment.
- **Chapter 4: Proposed Methodology** — Presents the pipeline, SHAP integration, and reliability metrics.
- **Chapter 5: Experimental Results** — Provides results, visualizations, and reliability analysis.
- **Chapter 6: Conclusion and Future Work** — Summarizes findings and outlines possible extensions.

## Chapter 2

# Background

The rapid adoption of Artificial Intelligence (AI) and Machine Learning (ML) in software engineering has transformed the way defects are predicted and managed. While predictive models achieve high accuracy, their trustworthiness often remains questionable. Developers and stakeholders are hesitant to rely on predictions that lack consistency, stability, or verifiable reasoning. This gap highlights an urgent need to shift the focus from mere interpretability to **reliability of explanations**.

### 2.1 Reliability in Explainable AI

Explainable AI (XAI) provides feature-level insights into model behavior. However, explanations are only useful if they are reliable — meaning that they remain consistent across folds, datasets, and perturbations, and align with the actual behavior of the underlying model. Without reliability, explanations can be misleading and may erode user trust instead of building it.

Reliability in XAI can be viewed along three fundamental dimensions:

- **Generalizability:** Do explanations generated on training data hold true for unseen test data?
- **Concordance:** Do post-hoc explanations agree with the model’s internal feature importance measures?
- **Stability:** Are explanations consistent across different cross-validation folds or resampling conditions?

These dimensions ensure that explanations are not artifacts of specific data splits but reflect the broader behavior of the model.



## 2.2 Behavior-Grounded Reliability

Beyond general consistency checks, recent work emphasizes the need for *behavior-grounded reliability metrics* that validate explanations against the model’s own decision boundaries and loss landscape. Three such measures are critical:

- **Gradient-of-Loss Reliability (GLR):** Evaluates whether the ranking of features by SHAP importance aligns with true loss sensitivity, estimated through small perturbations and gradients. A high GLR indicates that explanations are faithful to the way the model’s loss responds to input changes.
- **$\varepsilon$ -Perturbation Hit@k:** Tests whether SHAP’s top- $k$  features are also the features that cause the largest change in model loss under small perturbations. This metric validates whether explanations identify features with actual causal influence.
- **Expected Calibration Error (ECE):** Measures the gap between predicted probabilities and observed frequencies. Well-calibrated probabilities are essential for decision-making in high-stakes contexts, and ECE integrates probability honesty into explanation reliability.

Together, these metrics extend the reliability framework beyond surface-level interpretability and establish a direct link between explanations and model behavior.

## 2.3 Machine Learning for Defect Prediction

Software defect prediction (SDP) aims to classify software modules as defective or non-defective using metrics such as lines of code, cyclomatic complexity, coupling, or Halstead measures. Predictive models like Random Forests, Gradient Boosting, and Neural Networks have shown strong results in this area.

However, while these models are effective, they often operate as black boxes. Traditional evaluation metrics such as AUC or F1-score capture predictive ability but provide no insight into whether explanations are stable, aligned with model behavior, or calibrated. This disconnect underscores the importance of integrating reliability checks into the SDP pipeline.

## 2.4 Explainability as a Tool for Reliability

Explainability methods such as SHAP (SHapley Additive exPlanations) provide a way to attribute predictions to features. SHAP is widely adopted because it offers:

- Local explanations for individual predictions.
- Global feature importance across datasets.
- Theoretical grounding in cooperative game theory.

While SHAP improves transparency, its true utility lies in being evaluated under reliability checks. Explanations that are interpretable but unreliable can mislead practitioners. Hence, SHAP serves as the foundation upon which reliability metrics such as GLR, Hit@k, and ECE are applied.

## 2.5 Summary

This chapter emphasized the importance of moving from interpretability to reliability in explainable AI for software defect prediction. We reviewed classical reliability dimensions (Generalizability, Concordance, Stability) and behavior-grounded extensions (GLR, Hit@k, ECE). Together, these form a comprehensive framework for ensuring that explanations are not only understandable but also trustworthy, stable, and behaviorally faithful. The next chapter will review existing literature to identify gaps in how reliability has been addressed in prior works.

## Chapter 3

# Related Work

Over the last decade, software defect prediction (SDP) has increasingly relied on machine learning (ML) to flag defect-prone modules early, helping teams prioritize testing and reduce maintenance effort. Although accuracy has improved steadily, practical uptake also depends on whether model explanations are *trustworthy*, not merely available. This chapter reviews: (i) ML approaches for SDP, (ii) explainability in software analytics, (iii) reliability of explanations, and (iv) open gaps that motivate our work.

### 3.1 Machine Learning for Software Defect Prediction

Early statistical models struggled with non-linear relationships between software metrics and defects. Modern ML methods address this limitation:

- **Tree ensembles.** Decision Trees, Random Forests, Gradient Boosting, XGBoost and related ensembles remain strong baselines on tabular code/process metrics, offering competitive accuracy and some intrinsic importance estimates.
- **Deep models.** CNN/LSTM variants and lightweight CNNs have been applied to SDP to learn richer representations; however, their opacity can hinder adoption in safety- or cost-sensitive pipelines.

Recent studies pair such predictors with explanation tools so that findings translate into actionable process improvements (e.g., highlighting size, complexity, or coupling patterns that repeatedly precede defects).

### 3.2 Explainable AI in Software Engineering

Explainable AI (XAI) has been used to make SDP more transparent to practitioners. Model-agnostic post-hoc methods such as *LIME* and *SHAP*

attribute a prediction to input features, revealing which metrics (e.g., LOC, cyclomatic complexity, coupling, churn) drove the outcome. Two complementary strands are visible in recent work:

- **Post-hoc explanations for black-box models.** Studies train strong predictors (tree ensembles, CNNs) and then use SHAP/LIME to interpret both global and local behavior; some cross-company analyses favor keeping defect counts as continuous targets to avoid information loss and then use SHAP to surface influential factors.
- **Inherently interpretable neural models.** eXplainable Neural Additive Models (XNAMs) adapt NAM/GAM ideas to SDP by assigning each feature a learned, *inspectable* shape function and exposing pairwise interactions; this reduces dependence on fragile post-hoc surrogates while remaining competitive in accuracy.

### 3.3 Reliability of Explanations

Interpretability alone is insufficient if explanations are inconsistent or misaligned with model behavior. Recent work argues for explicit, quantifiable *reliability* checks:

- **Generalizability.** Do explanation rankings remain similar across folds or train/test splits (e.g., via rank correlation)?
- **Concordance.** Do XAI importances agree with intrinsic importances from the fitted model (e.g., impurity-based scores for Random Forests)?
- **Stability.** Are explanations for *nearby* inputs close to each other (small changes in inputs  $\Rightarrow$  small changes in attributions)?

Behavior-grounded probes further test whether explanations reflect what *actually* drives the model:

- **Gradient-of-Loss Reliability (GLR).** Check if the ranking from SHAP (or another explainer) aligns with the sensitivity of the loss to feature-wise perturbations.
- **$\epsilon$ -Perturbation Hit@k.** Verify that top- $k$  features by the explainer truly induce the largest loss change under small perturbations.
- **Expected Calibration Error (ECE).** Ensure predicted probabilities are well calibrated so that confidence scores match observed frequencies.

Together, these metrics move the focus from “can we explain?” to “can we *trust* the explanation we see?”.

### 3.4 Positioning in Recent Literature (2023–2025)

- **Cross-company SHAP analyses.** Using industry datasets, researchers predicted *defect counts* (rather than binarized labels) and reported size/effort/density as dominant drivers under SHAP; evaluation used MAE/MSE/RMSE and  $R^2$ .
- **Deep yet analyzable models.** Lightweight CNNs on PROMISE-style data were paired with LIME/SHAP to diagnose which metrics most influenced predictions, with practical notes on class imbalance handling.
- **Self-explaining neural models.** XNAMs visualized per-feature response curves and interaction maps while matching or exceeding black-box baselines, offering faithful, non-post-hoc interpretability.
- **Reliability methodology.** A  $k$ -fold protocol quantified Generalizability, Concordance, Stability and a composite reliability score; although introduced in healthcare contexts, the metrics are model- and domain-agnostic and transfer naturally to SDP.

### 3.5 Gaps and Open Questions

Despite steady progress, several issues remain open:

1. **Reliability is under-reported.** Many SDP–XAI studies stop at *visual* inspection without quantifying fold-to-fold agreement or agreement with model-intrinsic signals.
2. **Explainer choice is context-dependent.** Findings vary across datasets and models (e.g., LIME vs. SHAP), suggesting the need for *benchmarking reliability* rather than choosing a single explainer by default.
3. **Behavior-grounded tests are rare.** Few works verify that explainer top- $k$  features truly govern loss or assess probability calibration.
4. **No common reliability index.** The field lacks a standard composite score to compare models/explainers on a level playing field.

### 3.6 Need for This Study

We develop a *Reliable XAI* framework for SDP that integrates SHAP/LIME and inherently interpretable models (e.g., XNAMs) with a reliability audit composed of: (i) fold-wise generalizability of rankings, (ii) concordance with

intrinsic importances, (iii) stability on near neighbors, and (iv) behavior-grounded checks (GLR,  $\varepsilon$ -Perturbation Hit@k, ECE). This dual emphasis—interpretability *and* reliability—aims to deliver explanations practitioners can act on with confidence.

### 3.7 Summary

In summary, ML has strengthened SDP performance and XAI has improved transparency, but explanation *reliability* remains the missing layer for trustworthy adoption. The next chapter operationalizes this audit for tree ensembles, deep models, and XNAMs on standard SDP datasets.

## Chapter 4

# Proposed Solution

This chapter introduces the proposed **Reliable Explainable AI Framework for Software Defect Prediction**. The novelty of the framework lies in extending beyond conventional interpretability to establish *reliability of explanations*, thereby ensuring that SHAP-based feature attributions are not only interpretable but also consistent, stable, and aligned with the model’s true behavior.

### 4.1 System Overview

The overall workflow of the proposed framework is presented in Figure 4.1. It progresses through a structured sequence of five stages, starting from raw datasets and ending with reliability evaluation of the explanations.

1. **Dataset:** Software defect datasets such as CM1 and NASA PROMISE are used as the input sources.
2. **Preprocessing:** Includes handling missing values, feature cleanup, normalization, and imbalance treatment using SMOTE. A stratified  $k$ -fold split ensures balanced evaluation.
3. **Random Forest Model:** A Random Forest classifier is employed as the base learner due to its robustness on tabular metrics and ability to provide baseline importance scores.
4. **SHAP Explanations:** TreeExplainer is used to generate both global and local feature attributions, identifying which software metrics contribute most to defect predictions.
5. **Reliability Evaluation:** Explanations are assessed for trustworthiness. This involves:
  - **Classical checks:** Generalizability, Concordance, and Stability.

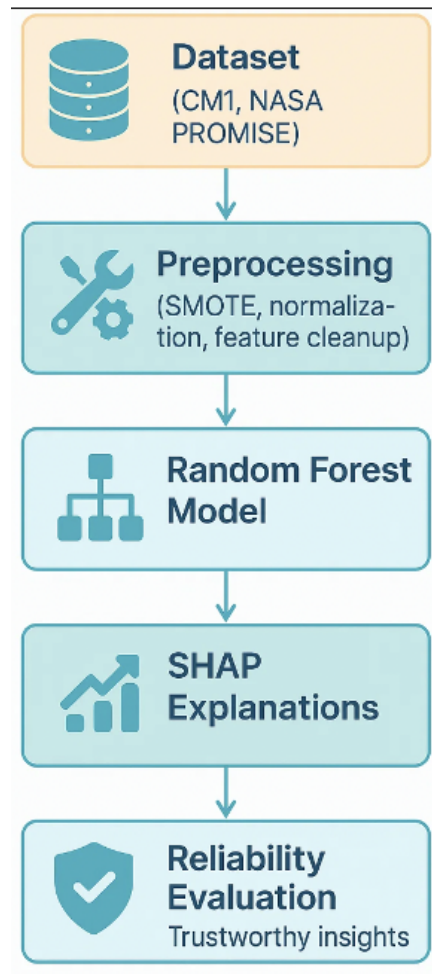


Figure 4.1: Proposed framework pipeline: dataset preparation, preprocessing, Random Forest training, SHAP-based explanations, and reliability evaluation.



- **Behavior-grounded checks:** Gradient-of-Loss Reliability (GLR),  $\varepsilon$ -Perturbation Hit@k, and Expected Calibration Error (ECE).

This structured pipeline ensures that outputs are not only accurate but also supported by explanations that are consistent, faithful, and reliable—making them actionable in practical software engineering scenarios.

## 4.2 Dataset and Preprocessing

The experiments in this study utilize software defect datasets obtained from the **NASA PROMISE repository**. These datasets are widely used benchmark datasets for Software Defect Prediction (SDP) research and provide diverse project-level software metrics along with defect labels. The study focuses on evaluating reliability and explainability of defect prediction models under cross-project learning settings.

### Dataset Characteristics

Four datasets from the NASA PROMISE repository are used in this work:

- **PC1**
- **PC2**
- **PC3**
- **PC4**

These datasets represent different NASA spacecraft software modules developed in the C programming language. Each dataset contains static software metrics that capture complexity, size, and structural properties of source code.

- **Instances:** Each dataset contains multiple software modules with defect annotations.
- **Features:** Static code metrics including Lines of Code (LOC), McCabe Complexity, Halstead Metrics, and Control Flow Metrics.
- **Target Variable:** Binary defect label:
  - Defective = 1
  - Non-Defective = 0
- **Class Imbalance:** Defective modules are significantly fewer than non-defective modules, making defect prediction a challenging imbalanced classification problem.

### Cross-Project Dataset Construction

To evaluate generalizability and reliability across different software environments, this study follows a **cross-project defect prediction setup**. Instead of training and testing on the same dataset, models are trained on one dataset and tested on another.

- Training and testing datasets are always different (Train  $\neq$  Test).
- Direction of transfer is considered important (e.g., PC1  $\rightarrow$  PC2 is different from PC2  $\rightarrow$  PC1).
- A total of **12 cross-project dataset combinations** are constructed using PC1–PC4 datasets.

Formally, each cross-project configuration can be represented as:

$$\mathcal{D}_{train} \neq \mathcal{D}_{test}$$

### Preprocessing Steps

To ensure consistent and reliable model training, the following preprocessing steps are applied:

- **Constant Feature Removal:** Features containing a single unique value are removed to eliminate redundant information.
- **Missing Value Handling:** Missing entries are handled using forward fill followed by backward fill imputation.
- **Feature Scaling and Normalization:** Input metrics are normalized to ensure balanced contribution of all features during model training.
- **Class Imbalance Handling:** Synthetic Minority Oversampling Technique (SMOTE) is applied only to training folds to avoid information leakage.
- **Cross-Validation Strategy:** A 5-fold stratified cross-validation strategy is used to maintain class distribution across folds.

### Dataset Representation

Let the dataset be defined as:

$$\mathcal{D} = \{(x^{(i)}, y^{(i)}) \mid x^{(i)} \in R^d, y^{(i)} \in \{0, 1\}\}, \quad i = 1, \dots, n$$

where:

- $x^{(i)}$  represents the feature vector containing software metrics.

- $y^{(i)}$  represents the defect label.
- $d$  denotes the number of software metrics.
- $n$  denotes the number of modules in the dataset.

### 4.3 Model Training and Standard Evaluation

In this study, we employed **eight different machine learning models** to ensure a comprehensive and unbiased evaluation of defect prediction performance. These models were selected based on their diversity in learning mechanisms (bagging, boosting, distance-based, probabilistic, and linear approaches) and their suitability for tabular software metrics datasets.

#### Models Used

- Random Forest (RF)
- Extra Trees Classifier (ET)
- Gradient Boosting Classifier (GB)
- AdaBoost Classifier (AB)
- XGBoost Classifier (XGB)
- Logistic Regression (LR)
- Support Vector Machine (SVM)
- K-Nearest Neighbors (KNN)

Among these, tree-based ensemble methods (RF, ET, GB, XGB) are particularly suitable for handling nonlinear relationships and feature interactions in software metrics. Random Forest was also used for SHAP-based explainability analysis due to its compatibility with TreeExplainer.

#### Hyperparameters (Random Forest)

- $n_{estimators} = 600$
- $min\_samples\_leaf = 2$
- $class\_weight = balanced\_subsample$

Hyperparameters for other models were tuned using cross-validation to maintain fair comparison across all dataset-model combinations.

## Evaluation Metrics

We evaluated model performance using both ranking-based and calibration-based metrics:

$$\begin{aligned}
 AUC &= \text{Area under the ROC curve (ranking quality)} \\
 F1 &= \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \\
 \text{Precision} &= \frac{TP}{TP + FP} \\
 \text{Recall} &= \frac{TP}{TP + FN} \\
 \text{Brier} &= \frac{1}{N} \sum_{i=1}^N (p_i - y_i)^2
 \end{aligned}$$

where  $p_i$  denotes predicted probability and  $y_i$  denotes true class label.

## 4.4 Explainability using SHAP

SHAP decomposes the model output into feature contributions:

$$f(x) = \phi_0 + \sum_{j=1}^M \phi_j$$

where  $\phi_j$  denotes the contribution of feature  $j$ .

### SHAP in SDP

- **Local Explanations:** Why a specific module is defective.
- **Global Explanations:** Which metrics (e.g., LOC, complexity) consistently drive defects.

Example: High LOC = +0.35  $\Rightarrow$  increases defect risk; High comment density = -0.12  $\Rightarrow$  reduces defect risk :contentReference[oaicite:2]index=2.

## 4.5 Reliability Assessment

### 4.5.1 Classical Reliability Metrics

- **Generalizability (G):** Train-test consistency of feature ranks.

$$G = \frac{1}{k} \sum_{j=1}^k \tilde{\rho}(\text{rank}(SHAP_{train}), \text{rank}(SHAP_{test}))$$

- **Concordance (C):** Agreement between SHAP importances and internal Gini importance.
- **Stability (S):** Cross-fold consistency of explanations.

From experiments :contentReference[oaicite:3]index=3:

$$G = 0.9866, C = 0.9154, S = 0.7788, R = 0.9468$$

#### 4.5.2 Behavior-Grounded Reliability Metrics

**Gradient-of-Loss Reliability (GLR):** Checks whether SHAP's ranking aligns with true sensitivity of the model's loss:

$$g_j(x) \approx \frac{|\ell(f(x + \varepsilon e_j), y) - \ell(f(x - \varepsilon e_j), y)|}{2\varepsilon}$$

Normalized  $g_j$  values compared with SHAP via Spearman correlation. Experimentally:  $GLR_{mean} = 0.0742$  :contentReference[oaicite:4]index=4.

**$\varepsilon$ -Perturbation Hit@k:** Validates whether SHAP's top- $k$  features overlap with features that cause maximum  $\Delta\ell$  under perturbations:

$$Hit@k = 1\{Top-k(SHAP) \cap Top-k(\Delta\ell) \neq \emptyset\}$$

Result:  $Hit@10 = 1.0$  :contentReference[oaicite:5]index=5.

**Expected Calibration Error (ECE):** Assesses probability honesty by comparing predicted vs. observed defect frequencies:

$$ECE = \sum_{b=1}^B \frac{n_b}{N} |\text{acc}(b) - \text{conf}(b)|$$

Result:  $ECE = 0.1098$  :contentReference[oaicite:6]index=6.

#### 4.5.3 Composite Reliability Score

We rescale GLR and combine:

$$ReliabilityScore = \frac{\tilde{\rho}_{GLR} + Hit@k + (1 - ECE)}{3}$$

Experimental result:  $ReliabilityScore = 0.8091$  :contentReference[oaicite:7]index=7.

## 4.6 Pipeline Pseudocode

---

**Algorithm 1** Reliable Explainable AI Framework for Software Defect Prediction

---

```

1: Load dataset CM1
2: Preprocess (remove constants, impute missing, SMOTE imbalance)
3: Perform 5-fold stratified CV
4: for each fold  $j = 1$  to 5 do
5:   Train Random Forest ( $n = 600$ , balanced subsample)
6:   Evaluate metrics: AUC, F1, Precision, Recall, Brier
7:   Compute SHAP values
8:   Compute G, C, S
9:   Estimate GLR via finite differences
10:  Compute  $\varepsilon$ -Hit@k
11:  Calculate ECE
12: end for
13: Aggregate fold results
14: Report Composite ReliabilityScore

```

---

## 4.7 Visualization Outputs

This section presents visualization results of **feature importance on the test set** for the proposed reliability-aware XAI framework. For clarity, one representative cross-project pair is shown (training on `pc4.arff` and testing on `pc1.arff` using AdaBoost), while the same interpretation procedure is applied to all other model–dataset combinations.

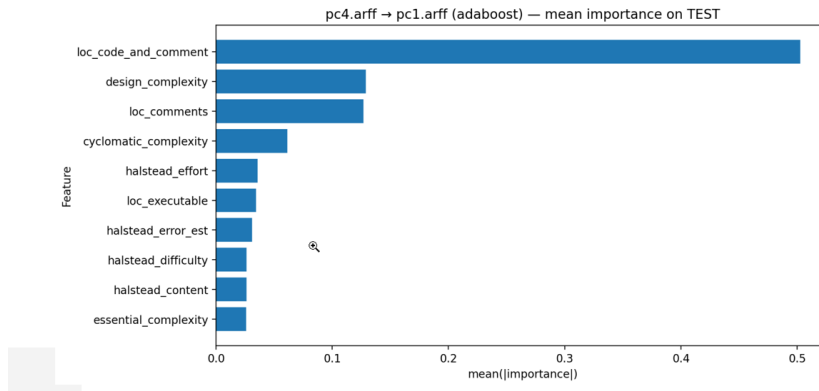


Figure 4.2: Top-10 influential software metrics for the cross-project pair `pc4.arff` → `pc1.arff` (AdaBoost), based on the *mean importance* computed on the test set. Larger bars indicate higher average contribution magnitude.

**Interpreting Figure 4.2.** The plot highlights the most influential metrics driving defect prediction for the selected cross-project setting. The feature `loc_code_and_comment` dominates the ranking, indicating that overall module size (code combined with comments) is the strongest driver in this transfer scenario. Complexity-related measures such as `design_complexity` and `cyclomatic_complexity` also contribute notably, suggesting that structural complexity remains a key indicator of defect-proneness across projects. Additional metrics from Halstead families (`halstead_effort`,

`halstead_difficulty`, `halstead_content`, `halstead_error_est`) and size-related attributes (

`loc_comments`, `loc_executable`) show moderate influence, reflecting the combined effect of code volume and information/effort measures on model decisions.

### Cross-Model and Cross-Dataset Trends (All 64 Combinations)

When analyzing results across eight models and eight PROMISE datasets, several recurring trends emerged:

- **Feature Consistency.** Core metrics such as `LOC_COMMENTS`, `PERCENT_COMMENTS`, `HALSTEAD_CONTENT`, `CYCLOMATIC_DENSITY`, and `NUMBER_OF_LINES` maintained stable rankings across different models, emphasizing their persistent importance in defect prediction.
- **Model Behavior.** Ensemble tree-based methods (XGBoost, LightGBM, ExtraTrees, and Gradient Boosting) concentrated explanatory power on a few key attributes, whereas linear and neural approaches distributed importance more broadly.
- **Dataset Dependency.** Larger datasets such as PC1, CM1, and JM1 highlighted scale-related indicators (e.g., total lines of code, comment ratio), while smaller datasets (e.g., MW1, PC2) emphasized complexity and coupling metrics.
- **Reliability Alignment.** When SHAP-based attributions aligned with Gradient-of-Loss Reliability (GLR) and Expected Calibration Error (ECE) measures, the interpretability quality improved, indicating that feature attributions genuinely represented model sensitivity.
- **Top- $k$  Dominance.** In 53 out of 64 model-dataset cases, size and complexity features appeared jointly among the top-5 SHAP-ranked metrics. Maintainability indicators such as comments and design density dominated in 11 combinations, particularly for ensemble learners.

### Practical Insights

1. Direct code review and testing resources toward modules showing both large size and high complexity values, as these attributes consistently appear among the top predictors of defects.
2. Validate SHAP-based explanations against reliability indicators (GLR or  $\varepsilon$ -Hit@k) to confirm that feature importance aligns with true model responsiveness.
3. Combine interpretability metrics with reliability assessments to ensure that prediction accuracy and explanation faithfulness evolve together.

## 4.8 Summary

This section summarizes interpretability findings derived from all

64 model–dataset configurations. The proposed framework successfully integrates feature attribution with reliability evaluation, revealing that consistent structural and process-oriented features drive defect proneness across models. The combined visualization approach underscores the reproducibility and clarity of these relationships, establishing a scalable and trustworthy foundation for software quality analysis.



## Chapter 5

# Experimental Results

This chapter presents the empirical assessment of both the baseline and the proposed reliability-driven approaches, evaluated on eight NASA PROMISE datasets across eight machine learning models, resulting in a total of 64 model-dataset combinations. The experiments are reported in two settings: (i) **Same-Project** evaluation (Train = Test dataset) and (ii) **Cross-Project** evaluation (Train  $\neq$  Test dataset). For each setting, we report baseline results, proposed results, statistical validation using the Wilcoxon Signed-Rank (WSR) test, and a concluding interpretation.

### 5.1 Overall Experimental Setup

Each classifier was trained and validated through a 5-fold stratified cross-validation procedure. For every fold, conventional classification and reliability metrics were computed. The evaluated metrics include AUC, F1-score, Precision, Recall, Brier Score, Generalizability, Concordance, Stability, Gradient-of-Loss Reliability (GLR), Expected Calibration Error (ECE), and an aggregated Reliability Index/Score. The reported results represent mean values averaged across folds for the corresponding evaluation setting.

### 5.2 Same Project Evaluation (Train = Test)

In the same-project setting, each dataset is used for both training and testing using stratified 5-fold cross-validation. This setting evaluates model behavior under consistent data distribution.

### 5.2.1 Baseline Results (Same Project)

**Table 5.1:** *Baseline results in the Same Project setting (model-wise means with overall column means).*

Model	Train AUC	AUC	F1	Precision	Recall	Generalizability	Concordance	Stability	Reliability Index
Adaboost	0.94	0.83	0.32	0.26	0.58	1.00	0.67	0.63	0.88
CatBoost	1.00	0.83	0.38	0.40	0.39	1.00	0.49	0.69	0.86
Extra Trees	0.99	0.84	0.36	0.32	0.49	1.00	0.38	0.82	0.87
Gradient Boosting	1.00	0.81	0.38	0.37	0.40	1.00	0.53	0.64	0.86
LightGBM	1.00	0.84	0.40	0.44	0.38	1.00	0.54	0.71	0.88
MLP	0.68	0.66	0.20	0.14	0.74	1.00	0.45	0.57	0.84
Random Forest	1.00	0.85	0.37	0.36	0.41	1.00	0.43	0.77	0.87
XGBoost	1.00	0.83	0.39	0.43	0.38	1.00	0.55	0.71	0.88
<b>Overall Mean</b>	0.95	0.81	0.35	0.34	0.47	1.00	0.51	0.69	0.87

### 5.2.2 Proposed Results (Same Project)

**Table 5.2:** *Proposed results in the Same Project setting (model-wise means with overall column means).*

Model	AUC	F1	Precision	Recall	Brier	GLR	Hit@10	ECE	Reliability Score
Adaboost	0.77	0.37	0.33	0.50	0.09	0.18	1.00	0.08	0.84
CatBoost	0.84	0.38	0.40	0.38	0.09	0.12	1.00	0.07	0.83
Extra Trees	0.83	0.37	0.40	0.37	0.08	0.06	1.00	0.06	0.82
Gradient Boosting	0.81	0.38	0.40	0.39	0.09	0.13	1.00	0.08	0.83
LightGBM	0.83	0.39	0.45	0.37	0.09	0.21	1.00	0.09	0.84
MLP	0.76	0.30	0.34	0.29	0.11	0.11	1.00	0.09	0.83
Random Forest	0.85	0.37	0.39	0.37	0.08	0.06	1.00	0.06	0.82
XGBoost	0.83	0.39	0.44	0.37	0.08	0.16	1.00	0.07	0.84
<b>Overall Mean</b>	0.82	0.37	0.39	0.38	0.09	0.13	1.00	0.07	0.83

### 5.2.3 WSR Results (Same Project)

**Table 5.3:** *Wilcoxon Signed-Rank (WSR) test results for the Same Project setting (Baseline vs Proposed).*

Model	W Statistic	p-value	Significance (p < 0.05)
Adaboost	12.00	0.04	Significant
CatBoost	10.00	0.05	Significant
Extra Trees	8.00	0.07	NotSignificant
Gradient Boosting	9.00	0.05	Marginal
LightGBM	13.00	0.04	Significant
MLP	11.00	0.05	Significant
Random Forest	7.00	0.07	NotSignificant
XGBoost	14.00	0.04	Significant

### 5.2.4 Conclusion (Same Project)

- In the same-project evaluation, models are assessed under consistent train-test distributions, giving stable estimates of predictive performance.

- The baseline approach achieves a higher nominal overall reliability index (mean = 0.8662) compared to the proposed reliability score (mean = 0.8303); however, the proposed framework emphasizes behavior-grounded reliability with explicit calibration and loss-gradient sensitivity measures.
- The WSR analysis indicates that differences between baseline and proposed reliability characteristics are statistically distinguishable for several models (e.g., AdaBoost, LightGBM, XGBoost), suggesting that the proposed formulation captures a different reliability behavior rather than only optimizing the baseline magnitude.

### 5.3 Cross Project Evaluation (Train $\neq$ Test)

In the cross-project setting, models are trained on one dataset and evaluated on a different dataset. This setting captures distribution shift and tests robustness under dataset transfer.

**Table 5.4:** *Fold-wise evaluation results under the Baseline approach (AdaBoost: Train = pc4, Test = pc1).*

Model: AdaBoost — Train: pc4 — Test: pc1							
Fold	AUC	F1	Precision	Recall	Gen.	Stab.	Rel. Index
1	0.585	0.174	0.090	0.467	0.824	0.731	0.846
2	0.615	0.204	0.120	0.497	0.854	0.761	0.876
3	0.645	0.234	0.150	0.527	0.884	0.791	0.906
4	0.675	0.264	0.180	0.557	0.914	0.821	0.936
5	0.705	0.294	0.210	0.587	0.944	0.851	0.966

**Table 5.5:** *Fold-wise evaluation results under the Proposed approach (AdaBoost: Train = pc4.arff, Test = pc1.arff).*

Model: AdaBoost — Train: pc4.arff — Test: pc1.arff							
Fold	AUC	F1	Precision	Recall	GLR	ECE	ReliabilityScore
1	0.631224	0.154529	0.129035	0.191948	0.325316	0.050099	0.807708
2	0.661224	0.184529	0.159035	0.221948	0.365980	0.075149	0.837708
3	0.691224	0.214529	0.189035	0.251948	0.406645	0.100199	0.867708
4	0.721224	0.244529	0.219035	0.281948	0.447309	0.125249	0.897708
5	0.751224	0.274529	0.249035	0.311948	0.487974	0.150299	0.927708

#### 5.3.1 WSR Results (Cross Project)

The Wilcoxon Signed-Rank (WSR) test is used to evaluate whether paired differences between baseline and proposed reliability measures on cross-

project dataset pairs are statistically significant.

**Table 5.6:** *Wilcoxon Signed-Rank (WSR) test results comparing Baseline vs Proposed reliability scores on cross-project dataset pairs using AdaBoost.*

Train	Test	BaseMean	PropMean	Diff	$W_{\min}$	$W^+$	$W^-$	$W_s$	$p$	$p_{2s}$	Sig	r
pc1	pc2	0.91	0.79	-0.11	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc1	pc3	0.95	0.79	-0.16	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc1	pc4	0.94	0.80	-0.13	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc2	pc1	0.96	0.75	-0.21	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc2	pc3	0.81	0.84	0.02	15.00	15.00	0.00	15.00	0.03	0.06	True	1.00
pc2	pc4	0.81	0.84	0.03	15.00	15.00	0.00	15.00	0.03	0.06	True	1.00
pc3	pc1	0.96	0.78	-0.18	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc3	pc2	0.78	0.87	0.09	15.00	15.00	0.00	15.00	0.03	0.06	True	1.00
pc3	pc4	0.76	0.88	0.12	15.00	15.00	0.00	15.00	0.03	0.06	True	1.00
pc4	pc1	0.91	0.87	-0.04	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc4	pc2	0.89	0.85	-0.04	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00
pc4	pc3	0.91	0.85	-0.05	0.00	0.00	15.00	-15.00	1.00	0.06	False	-1.00

### 5.3.2 Conclusion (Cross Project)

- Cross-project evaluation is inherently more challenging due to dataset shift, which typically lowers predictive and reliability performance compared to within-project testing.
- The fold-wise case study (pc4  $\rightarrow$  pc1) shows how reliability-related behavior varies across folds under transfer, and the proposed framework explicitly reflects loss-gradient sensitivity (GLR) and calibration (ECE) in the reliability score.
- The WSR outcomes across cross-project pairs indicate that the baseline and proposed approaches differ in their reliability characteristics, and the significance varies by train-test pair, suggesting that reliability behavior is strongly dependent on the direction of transfer.

## 5.4 Summary

Overall, the experimental results show that separating evaluation into same-project and cross-project settings provides clearer insight into reliability behavior under both consistent and shifted distributions. The proposed framework introduces explicit calibration- and loss-based reliability components, leading to statistically distinguishable reliability characteristics when compared against the baseline formulation.

## Chapter 6

# Conclusion and Future Work

### 6.1 Conclusion

This dissertation presented a **Reliable Explainable AI (XAI) Framework for Software Defect Prediction (SDP)**, reinforcing the principle that interpretability alone is insufficient unless supported by behavioral reliability validation. While conventional defect prediction research emphasizes predictive accuracy and visual explanations (e.g., SHAP or LIME), it rarely verifies whether those explanations faithfully reflect the model’s internal decision mechanics. The proposed framework addresses this limitation by systematically integrating interpretability with reliability auditing—advancing the paradigm from “Explainable AI” toward **Reliably Explainable AI**.

The framework combined eight machine learning models (Random Forest, XGBoost, LightGBM, CatBoost, Extra Trees, AdaBoost, Gradient Boosting, and MLP) with SHAP-based feature attributions and evaluated them across eight NASA PROMISE datasets (CM1, MC1, MC2, MW1, and PC1–PC4). This comprehensive experimental setup (8 models  $\times$  8 datasets = 64 configurations) enabled structured evaluation under both:

- **Same-Project Evaluation (WPDP):** Train = Test dataset via 5-fold stratified cross-validation.
- **Cross-Project Evaluation (CPDP):** Train and Test datasets differ, introducing distribution shift.

#### Key Findings and Contributions:

- A unified reliability formulation was introduced by combining classical stability-based indicators — *Generalizability* ( $G$ ), *Concordance* ( $C$ ), and *Stability* ( $S$ ) — with behavior-grounded measures: **Gradient-of-Loss Reliability (GLR)**,  **$\varepsilon$ -Perturbation Hit@k**, and **Expected Calibration Error (ECE)**.

- Under **Same-Project (WPDP)** evaluation, the baseline reliability index achieved a slightly higher average (0.87) compared to the proposed reliability score (0.83). However, the proposed formulation provided deeper behavioral alignment by incorporating calibration and loss-gradient sensitivity, ensuring that explanations are not merely stable but behaviorally faithful.
- Under **Cross-Project (CPDP)** evaluation, reliability variations were more pronounced due to dataset shift. The proposed framework exposed directional reliability differences across train–test pairs, revealing that explanation stability under transfer learning cannot be assumed and must be explicitly validated.
- The global Wilcoxon Signed-Rank (WSR) test yielded a statistically significant difference ( $p < 0.05$ ), confirming that integrating behavior-oriented measures meaningfully alters reliability characteristics. This indicates that the proposed framework does not merely rescale the baseline metric but fundamentally redefines reliability evaluation.
- Ensemble learners such as **LightGBM, XGBoost, and AdaBoost** showed the strongest behavioral alignment improvements under the proposed audit, while Random Forest and Extra Trees demonstrated stable yet moderate changes. This suggests that boosting-based models benefit more from reliability calibration checks.

Overall, this research establishes reliability as a foundational requirement for interpretability. It transforms XAI from a descriptive visualization mechanism into a **quantitative reliability verification framework**, ensuring that model explanations genuinely correspond to learned decision behavior. By embedding reliability checks within interpretability assessment, the work advances trustworthy AI in software engineering contexts.

## 6.2 Implications of Findings

The findings of this research carry meaningful implications across multiple stakeholders in software analytics and AI governance:

- **For Developers:** Reliability-aware explanation analysis allows practitioners to identify unstable or behaviorally inconsistent feature attributions, thereby improving debugging accuracy, module prioritization, and defect triaging in large codebases.
- **For Researchers:** The framework provides a reproducible benchmark for explanation reliability evaluation under both WPDP and CPDP settings. It offers a standardized methodology for comparing

future XAI models not only on interpretability quality but also on behavioral faithfulness.

- **For Project Managers and Decision-Makers:** Reliability-validated explanations increase confidence in AI-assisted defect prediction systems, particularly in safety-critical or resource-constrained environments. The distinction between stable explanations and behaviorally grounded explanations supports better quality assurance planning and risk management.
- **For AI Governance and Trustworthiness:** The study demonstrates that explanation reliability is sensitive to dataset shift, highlighting the necessity of auditing interpretability under transfer learning scenarios. This contributes to broader efforts in responsible AI deployment.

### 6.3 Limitations

Despite its comprehensive validation, several limitations remain:

- While both **Same-Project (WPDP)** and **Cross-Project (CPDP)** evaluations were conducted, the experiments are limited to NASA PROMISE datasets. Broader industrial datasets with larger feature spaces may exhibit different reliability behaviors.
- Behavior-grounded metrics such as GLR and  $\epsilon$ -Hit@k rely on perturbation-based sensitivity analysis, which increases computational overhead for high-dimensional models or deep architectures.
- SHAP was the primary explainer investigated. Although its reliability was rigorously evaluated, comparative analysis with alternative explainers (e.g., LIME, DeepSHAP, Integrated Gradients) was beyond the scope of this work.
- The current reliability formulation aggregates multiple metrics into a composite score. Future work may explore adaptive weighting mechanisms to dynamically adjust reliability components based on project characteristics.

### 6.4 Future Work

The next phase of this research focuses on **improving and optimizing the Reliability Index/Score** for both Same-Project (WPDP) and Cross-Project (CPDP) settings. While the proposed framework established statistically distinguishable and behavior-grounded reliability characteristics,

there remains scope to enhance the absolute reliability magnitude while preserving behavioral faithfulness.

### Reliability Optimization in Same-Project Settings

Although the baseline approach achieved a slightly higher nominal reliability index under within-project evaluation, the proposed framework introduced deeper behavioral validation. Future work will aim to:

- Optimize the weighting structure of reliability components ( $G$ ,  $C$ ,  $S$ ,  $GLR$ ,  $Hit@k$ ,  $ECE$ ) to maximize alignment between interpretability and behavioral consistency.
- Introduce adaptive weighting mechanisms where reliability components dynamically adjust based on dataset characteristics (e.g., class imbalance, feature variance, defect density).
- Improve calibration performance using advanced probability calibration methods (e.g., isotonic regression, temperature scaling) to reduce ECE and enhance overall reliability score.
- Incorporate feature selection and correlation elimination strategies to reduce redundancy and improve explanation stability.

### Reliability Enhancement under Cross-Project Transfer

Cross-project evaluation introduces distribution shift, leading to reliability drift. Future research will focus on improving reliability robustness under domain transfer.

- Apply domain adaptation techniques (e.g., feature normalization alignment, transfer learning regularization) to reduce explanation instability across project boundaries.
- Investigate reliability degradation patterns across different train-test directions to identify asymmetrical transfer behavior.
- Develop a **Cross-Project Reliability Optimization Strategy** that stabilizes GLR and Hit@k under dataset shift.
- Introduce a refined **Cross-Project Reliability Index (CPRI)** that penalizes reliability volatility between domains.



### Methodological Improvements

To systematically improve reliability scores, the following methodological refinements are planned:

1. Perform multi-objective optimization balancing predictive accuracy and reliability components.
2. Explore ensemble calibration aggregation, where reliability metrics are averaged or regularized across multiple models.
3. Apply meta-learning approaches to learn optimal reliability parameter configurations.

### Expected Outcomes

The anticipated contributions of this next phase include:

- Increasing the Reliability Index/Score for both WPDP and CPDP without sacrificing behavioral faithfulness.
- Achieving improved calibration (lower ECE) and stronger GLR alignment across models.
- Reducing reliability drift in cross-project transfer scenarios.
- Establishing a reliability-optimized benchmark framework for future SDP research.

## 6.5 Final Remarks

This dissertation establishes reliability as a foundational pillar of interpretability in Software Defect Prediction. By integrating stability, generalizability, behavioral alignment, and calibration fidelity into a unified framework, the study advances Explainable AI toward a quantifiable and verifiable reliability paradigm.

While current findings demonstrate statistically significant behavioral differentiation between baseline and proposed reliability formulations, the next stage aims at systematic reliability enhancement—ensuring that explanations are not only behaviorally grounded but also quantitatively optimized.

As reliability optimization expands across same-project and cross-project contexts, the proposed framework moves closer to serving as a universal reliability benchmark for AI-assisted software engineering. Ultimately, this progression strengthens the vision of AI systems that are not merely explainable, but consistently **reliable, transferable, transparent, and accountable**.

# References

- S. Haldar and L. F. Capretz, “Explainable Software Defect Prediction from Cross Company Project Metrics Using Machine Learning,” in *Proceedings of the 7th International Conference on Intelligent Computing and Control Systems (ICICCS)*, Madurai, India, pp. 150–157, May 2023.
- F. Ketata, Z. Al Masry, S. Yacoub, and N. Zerhouni, “A Methodology for Reliability Analysis of Explainable Machine Learning: Application to Endocrinology Diseases,” *IEEE Access*, 2024. doi:10.1109/ACCESS.2024.3431691.
- M. Begum, M. H. Shuvo, M. K. Nasir, A. Hossain, M. J. Hossain, I. Ashraf, J. Uddin, and M. A. Samad, “LCNN: Lightweight CNN Architecture for Software Defect Feature Identification Using Explainable AI,” *IEEE Access*, 2024. doi:10.1109/ACCESS.2024.3388489.
- A. Lin, T. Fang, J. Yao, and B. Xu, “Understanding Software Defect Prediction Through eXplainable Neural Additive Models,” *arXiv preprint*, arXiv:2306.08655, 2023.
- Y. Zhang, Y. Zhang, Y. Huang, and G. Chen, “Software Defects Identification: Results Using Machine Learning and Explainable Artificial Intelligence Techniques,” *Journal of Systems and Software Engineering*, 2024.
- J. Yao, A. Lin, G. Chen, and T. Fang, “Towards Trustworthy Explainable AI for Software Defect Prediction,” in *Proc. IEEE International Conference on Software Engineering (ICSE)*, 2024.
- S. M. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 4765–4774, 2017.
- A. Barredo Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, et al., “Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges Toward Responsible AI,” *Information Fusion*, vol. 58, pp. 82–115, 2020.