

A Unified Deep Learning Framework for Software Bug Category Prediction Using Mixed Embeddings

Veena Kulkarni ^{1*}, Dr. Anand Khandare ²

¹ Thakur College of Engineering and Technology, Kandivali (E), Mumbai India. veena.kulkarni@thakureducation.org

² Thakur College of Engineering and Technology, Kandivali (E), Mumbai , India.

Abstract: Software bug classification has been considered an important aspect on enhancing software maintenance and the defect management process within large rescue software development projects. Bug reports may include a wide range of information e.g. textual description, code snippets and stack trace information which can give invaluable information in determining the nature of the software defects. But the main limitation of the traditional methods of classification of bugs is based on the analysis of textual characteristics, and this approach does not allow them to capture the structural information that can be found in the code and execution traces. The paper is a proposal of a consolidated deep learning architecture to predict the category of bugs in software based on mixed set of embeddings of various elements of bug reports. The proposed implementation involves the use of BERT that creates contextual embeddings based on natural language descriptions and CodeBERT that extracts embeddings on a code snippet and stack trace information. Such embeddings are summed in a mixed embedding model, which is subsequently utilized to prepare a classifier to act classifying bugs in a category of UI/UX, a category of performance and a category of security-related bugs. Evaluation of open source bug data set of size 6151 bug reports were performed in an experimental manner. The findings indicate that the proposed model has a better classification performance as compared to the baseline model based on GloVe embeddings and LSTM. These additions containing textual, code, and stack trace embeddings make feature representation much more useful and enhance the usability of automated bug-triaging systems in the software development setting..

Keywords: Software Bug Categorization Mixed Embeddings Deep Learning BERT CodeBERT Bug Report Classification Software Defect Prediction Bug Triaging.

1. Introduction

The rapid growth of modern software systems, especially large-scale distributed, web-based, and cloud-enabled applications, has significantly increased the complexity of software development and maintenance. Modern software systems consist of multiple interacting modules, third-party libraries, APIs, and cross-platform interfaces that must function together seamlessly. As a result, the likelihood of software defects occurring during development and deployment has increased considerably. Software defects or bugs are unavoidable in software engineering and may arise at any stage of the software development life cycle, including design, coding, integration, testing, and maintenance. Managing these bugs efficiently is therefore a crucial aspect of ensuring software reliability, security, and overall system performance.

With the widespread adoption of collaborative development platforms such as GitHub, GitLab, and Bitbucket, large-scale software projects generate thousands of bug reports daily. These bug reports are typically stored in bug tracking systems such as Bugzilla, Jira, and Redmine. The reports usually contain textual descriptions, error logs, screenshots, stack traces, system configurations, and user comments describing the encountered issues. Because bug reports are primarily written in natural language, extracting meaningful information from them becomes challenging. Consequently, the manual analysis and classification of these bug reports by developers and testers can become time-consuming and error-prone, especially in projects with large user bases and continuous updates.



Due to the programming errors, due to inadequate testing, requirements are wrong or incomplete, invalid access to users, User Interface colour etc. There are a vast number of software bugs reported on a daily basis on big projects by the testers/users. These bug reports mainly consist of bug descriptions in text [4]. The bugs are of varied nature in terms of their occurrence and impact. Bug categorization is done manually by developers or testers, who would: Analyse the bug report,

identify affected components, assign a severity level, determine a bug category (e.g., UI, Performance, Security). This is a very time-consuming process; it can delay the performance of the project. The Bug Life Cycle ideally consists of the following:

1. Identification of bugs: Bugs are identified through testing, user feedback.
2. Logging in any bug management system: Developers log the bug with detailed information about its type, impact and occurrence.
3. Assignment to developers: Bugs are assigned to developers for fixing.
4. Fixing of bugs: Developers resolve the bug by modifying the code.
5. Testing of bug fixes: The fixed code undergoes testing to ensure that the bug is truly fixed.
6. Closure: Once the bug is verified to be solved, the bug is closed.

Software bug categorization is a crucial part of software testing that involves identifying the type of bug while logging in the bug management system. As software systems grow in complexity, automated bug classification/categorization into different prescribed bug types as per the projects helps developers to prioritize fixes. This helps to improve overall software bug triaging [5]. Automated bug categorization not only accelerates the debugging process but also helps organizations track recurring issues, and also allocate resources efficient. The various types of bus is given in Table 1.

Table 1 Bug types

Bug Type	Description	Examples	Impact on System
UI/UX Bugs	Issues related to the graphical user interface, layout design, usability, or navigation of the software application that affect the user experience.	Misaligned buttons, incorrect color scheme, broken links, overlapping text, unresponsive UI elements.	Reduces usability and user satisfaction, may lead to confusion or incorrect system usage.
Performance Bugs	Defects that affect the efficiency, speed, and resource utilization of the system during execution. These bugs often arise due to inefficient algorithms, memory leaks, or improper resource management.	Slow page loading, high CPU usage, excessive memory consumption, delayed response time, system lag under heavy workload.	Decreases system efficiency, increases resource consumption, and may cause system slowdown or crashes.
Security Bugs	Vulnerabilities or weaknesses in the system that may allow unauthorized access, data manipulation, or malicious attacks.	SQL injection, authentication bypass, insecure API calls, weak encryption, improper access control.	May lead to data breaches, privacy violations, system compromise, or unauthorized system access.

The main contributions of this research are summarized as follows. First, a comprehensive literature review of existing studies on software bug categorization has been carried out in order to understand the current techniques, challenges, and research trends in automated bug classification. Second, a deep learning-based framework is proposed for software bug categorization. In the proposed approach, BERT is used to generate textual embeddings from the bug descriptions, while CodeBERT is utilized to extract embeddings from code snippets and stack trace information contained in the bug reports. The embeddings obtained from these two sources are concatenated to form a unified feature representation, which is then provided as input to a BERT-based classifier to predict the appropriate bug category. Third, for comparative analysis, the GloVe model is applied to generate non-contextual word embeddings from bug descriptions, and these embeddings are used as input to an LSTM classifier [7]. The performance of this baseline approach is then compared with the proposed mixed-embedding deep learning framework to evaluate its effectiveness.

The remainder of this paper is organized as follows. Section 2 presents the related work on software bug categorization and previously proposed machine learning approaches. Section 3 describes the

proposed framework for software bug category prediction in detail. Section 4 presents the experimental results and discussion, including a comparative analysis between the proposed method and existing approaches. Finally, Section 5 provides the conclusion of the study and outlines potential directions for future research.

2. Related Work

The classification of software bugs and prediction of their presence have been the focus of studies over the last few years, notably as the machine learning and deep learning approaches continue to grow in use in software engineering. Various studies have investigated various methods of bug report analysis and automatic categorization into some set of predetermined categories.

Meng F., Huang R., and Wang J. [1] provided a thorough review of the research about software defects in the last three years, relying on deep learning. Their work sums up different deep learning systems in the software defect prediction and identification. Models that were discussed by the authors include **TextCNN** and TextRNN to predict software defects whereas models such as LSTM, BiLSTM, DNCC and DAKSM were used to identify defects. Moreover, an **Atten-CRNN** model was employed in automated assignment of bugs. The research paper also conducted a comparative study of these models against conventional methods and indicated the increased accuracy of deep learning-based methods.

Rizal Broer Bahaweres et al. [7] suggested a bug classification framework that integrates an implementation of **Long Short-Term Memory (LSTM)** networks with a word embedding strategy. Their study data was obtained on some of the Java based open source software projects. Word embeddings facilitated the model to encode semantic associations in the description of bugs, which enhanced the performance of the LSTM model in classification.

An empirical study conducted by Gemma Catolino et al. [2] on the 1,280 reports of software bugs on 119 open software projects such as Mozilla, Apache and Eclipse was empirical. The analysis was aimed at determining and classifying the various reported software bugs. According to their analysis, the authors discovered a total of nine different types of bugs and formulated a classification model that is used to classify the bug reports automatically under the nine different types of bugs. The obtained F1-score of the proposed model was 64% and the value of AUC-ROC was 74%, which indicates that automated bug classification without human involvement is possible with the help of machine learning.

The study conducted by Nevendra, Meetesh Singh, and Pradeep [3] presents the research focused on whether deep learning methods are effective in predicting software defects or not. Various deep learning models were utilized in a number of datasets in their work to compare their strengths and weaknesses. The authors have studied the problem of predicting software defects using different machine learning methods in depth, and their results show that deep learning methods tend to outperform other methods based on traditional machine learning since they are able to learn intricate patterns of software defects in both textual and formatted data automatically.

The study by Tan, L., Liu, C., Li, Z., et al. [6] analyzed the bug reports in Linux operating system repositories and classified the bugs into three predominant groups including: **memory bugs, semantic or security bugs, and concurrency bugs**. The authors also examined the correlation between the types of bugs and the causes of bugs. Their results showed that a very large percentage of software crashes had to do with bugs that concerned the memory. Moreover, the research examined the connection between the root causes of bugs and the level or the severity of their effect on the systems.

A proposed model suggested by Jyoti Prakash Meher et al. [8] is a bug classification model based on deep learning and includes eight classes of software bugs. It also included the generation of a big annotated dataset, in which annotation was carried out based on the Earth Mover Distance (EMD) method which determines the relevant keywords in each category of bugs. The authors used some variants of attention-based transformer models, i.e. **Transformer, CodeBERT, BERT, and DistilBERT**, to perform automated bug classification. The proposed method attained a high F1-score of 84.74% that illustrates that transformer-based structures are effective when dealing with bug categorization tasks.

Koksal O. et al. [9] had a very comprehensive background of the usage of machine learning and deep learning techniques in software bug classification. The paper has talked about the different algorithms, methods of embedding words and transformer-based prebuilt models like **BERT**. The authors compared the analyses of 19

various research works and carried out a comparative analysis on the basis of the datasets, classification algorithms, and methods of extracting features applied. Their results have indicated the performance of classification using the F1-scores of between 62% and 88% in the first indication of effectiveness of the latest NLP-based methods in software defect analysis.

Yadav et al. [10] examined the problem of and the prospects of important defects software of the dataset. The authors emphasized that the datasets with a small number of incidences of some categories of bugs might result in false classification and high costs of prediction. They examined two major methods of dealing with this problem: dataset-level techniques and algorithm-based techniques. According to their findings, the ensemble learning techniques mixed with data level strategies are an effective way of addressing the issue of class imbalance.

The issue of class imbalance was also considered in the problem of the classification of the severity of bugs with multi-classes by N. K. S. Roy et al. [11]. Data of 17 software projects were utilized in their experiments, which investigated the usefulness of various techniques of feature extraction programs. The authors examined unigrams, bigrams and feature selection approaches used together with Support Vector Machine (SVM) classifiers. Their findings revealed that the performance of bigram and feature selection have a strong impact on the performance of classification in relation to a baseline classifier.

The researchers created a better variant of BERT model; namely, Yinhan Liu et al. [12] proposed RoBERTa that can be used to study the problems of language representation. The model was trained using a much larger corpus and tested using several benchmark sets. RoBERTa also uses a dynamic masking strategy wherein the masking pattern is varied during training, and which enables better contextual representation learning. The model was shown to be performing better than BERT and other state-of-the-art models on benchmarking tasks including the GLUE, RACE and the SQuAD.

H. A. Ahmed et al. [13] suggested a categorization approach of software bugs which can be done in an automatic fashion through supervised learning. They used a manual analysis of the bug repository as part of their application to produce a dataset of classification models. The extraction of features was conducted through the application of TF-IDF technique and the use of SMOTE algorithm was used to overcome the problem of class imbalance. To arrive at a conclusion on which type of the classifiers would work best in classifying bugs, the authors tested various classifiers, which include; **Naiv Bayes, Decision Tree, Resurrent Forest, and Logistic Regression.

Table 2: Summary of Existing Bug Classification Techniques and Their Limitations

Citation	Technique / Model	Embedding / Feature Extraction	Dataset Used	Performance Reported	Limitations
[14]	Deep Learning Models (TextCNN, TextRNN, LSTM)	Word Embedding	Public defect datasets	Improved accuracy compared to traditional ML models	Limited contextual understanding of bug descriptions
[15]	Machine Learning Classification	TF-IDF	Mozilla, Apache, Eclipse bug reports	F1-score: 64%, AUC-ROC: 74%	Moderate classification accuracy and limited semantic representation
[16]	Deep Learning Comparative Study	Feature-based representations	Multiple defect datasets	DL models outperform ML methods	No unified framework for bug category prediction
[17]	Statistical Bug Analysis	Manual feature extraction	Linux OS bug dataset	Identification of memory and concurrency bugs	Focus on analysis rather than automated classification
[18]	LSTM-based Bug Classification	Word2Vec / Word Embedding	Java open-source repositories	Improved textual bug classification	Non-contextual embeddings reduce semantic understanding

[19]	Transformer-based Models (BERT, CodeBERT, DistilBERT)	Contextual Embeddings	Annotated bug dataset	F1-score: 84.74%	Requires large datasets and high computational resources
[20]	ML/DL Comparative Study	Various embeddings and NLP features	Multiple software bug datasets	F1-score between 62–88%	Survey-based study without proposing new framework
[21]	Supervised Learning (NB, DT, RF, LR)	TF-IDF with SMOTE	Constructed bug repository dataset	Moderate classification accuracy	Limited ability to capture contextual and semantic relationships

Therefore, the literature review demonstrates the increasing popularity of machine learning, natural language processing, and deep learning to classify bugs automatically. Although older machine learning models have shown moderate results, recent research indicates that deep learning model and transformer structures should deliver better results in terms of accuracy since they can learn semantic relationships in bug reports. Nevertheless, issues like unequal measure of data, textual description, and features representation are also unresolved research problems. The obstacles encourage researchers to come up with more robust systems that integrate various embedding methods and deep learning models to make better predictions of software bugs categories.

3. Proposed Framework

A collective scheme of predicting software bugs is offered on the basis of the discussion of the available literature and the shortcomings found in the research conducted earlier. Most of the current methods mainly emphasize the textual representation of the bug report and disregard other aspects of the information like the code snippets and stack traces. Nevertheless, the bug reports can include several sources of information that can be used as valuable information to categorize the bug appropriately.

The proposed framework would make use of all the elements of the bug report description such as textual data, pieces of code, and information about stack traces. Bug text description is done with the help of the BERT model which creates textual embeddings of the context which reflects the semantic meaning of the reported bug problem. Likewise, CodeBERT can be used to extract the embeddings of code snippets and stack trace data because it is particularly created to comprehend the programming language structure and code semantics. These various sources are then used to generate embeddings which are subsequently summed together to form one sense of representation. Embeddings This symbiotic representation is presented as an input to a deep learning classifier that predicts the type of bug to enhance the accuracy and the robustness of the classification process. The proposed methodology has been provided in a diagrammatical form in Figure 1..

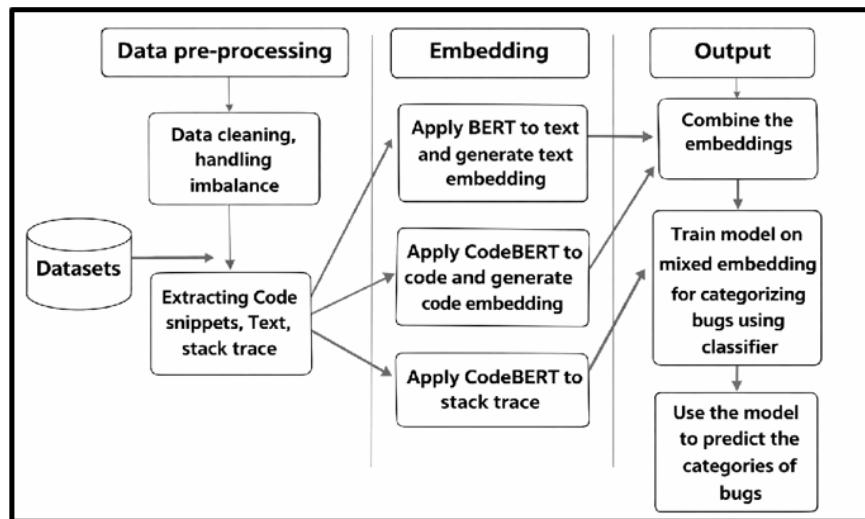


Figure 1. Proposed Architecture

3.1 Dataset used

The bug data employed in the research was sourced via an open-source bug-tracking system, Bugzilla, which is a collection of bug reports of popular open-source software projects. This database is composed of 5,651 bug reports to test and 500 bug reports to train. As a complete set of labels that could be used in the training step was not immediately available, the training bug reports were first labeled with the help of a large language model (LLM) and then manually maintained to provide the same consistency and reliability as the original open-source dataset.

The data is handled in form of CSV (Comma-Separated Values) and it includes a set of attributes with each bug report. These attributes have fields like title, description, date, severity and priority and other metadata on reported issue. These qualities give contextual valuable information which helps in comprehending the nature and effects of the bugs.

The data has both brief descriptions and lengthy descriptions of the reported bugs. In the suggested method, textual fields that are the long description and short description [16] are integrated into one textual summary. This brings on board much additional contextual information regarding the bug report and serves to key in as the main textual input to the stages of embedding and classifying bugs in the proposed software bug categorization framework.

Table 3: Distribution of Bug Reports Across Categories

Bug Category	Number of Bug Reports	Percentage (%)
UI/UX Bugs	2,145	34.87%
Performance Bugs	2,012	32.71%
Security Bugs	1,994	32.42%
Total	6,151	100%

The distribution of the bug reports in the three categories that are taken into account in this paper is given in Table 3. The bug report data includes the UX/UI bug reports, performance-related bug reports, and data breaches. The number of bug reports per category is relatively equal and this assists in minimising bias in model training and evaluation. Balanced data set is used to confirm that the suggested deep learning model will be capable of learning the patterns connecting each type of bugs and enhance the credibility of the results of classification.

3.2 Methodology

The proposed methodology focuses on software bug category prediction using mixed embeddings generated from textual descriptions, code snippets, and stack trace information present in bug reports. The overall process consists of multiple stages including data preprocessing, feature extraction, embedding generation, embedding fusion, model training, and performance evaluation. The objective of the methodology is to combine semantic information from natural language and programming artifacts to improve the classification accuracy of bug categories.

3.2.1 Data Preprocessing

The first stage of the methodology involves preprocessing the bug dataset to remove noise and normalize the textual information. Bug descriptions collected from open-source repositories often contain HTML tags, URLs, numerical values, and unnecessary formatting elements that may affect the quality of feature extraction. Therefore, several Natural Language Processing (NLP) preprocessing steps are applied in the methodology as shown in Figure 2.

Initially, the bug descriptions are converted into lowercase format and HTML tags are removed. Tokenization is then performed using the NLTK library, which splits the bug description into individual tokens or words. Stop words such as “the”, “is”, and “and” are removed as they do not contribute significantly to semantic meaning. Additionally, URLs and numerical values are eliminated from the text. Stemming and lemmatization techniques are applied to reduce words to their root forms, thereby improving feature consistency.

The cleaned textual representation of the bug report can be mathematically represented as

$$T_{clean} = f_{pre}(Tra)$$

Where $Tra()$ represents the original bug description text and $f_{pre}(\cdot)$ represents the preprocessing function including tokenization, stop word removal, and normalization.

This transformation is expressed in Equation (1).

$$T_{clean} = f_{pre}(Tra()) \quad (1)$$

Equation (1) shows that the preprocessing function transforms the raw textual description into a cleaned and normalized representation suitable for further processing.

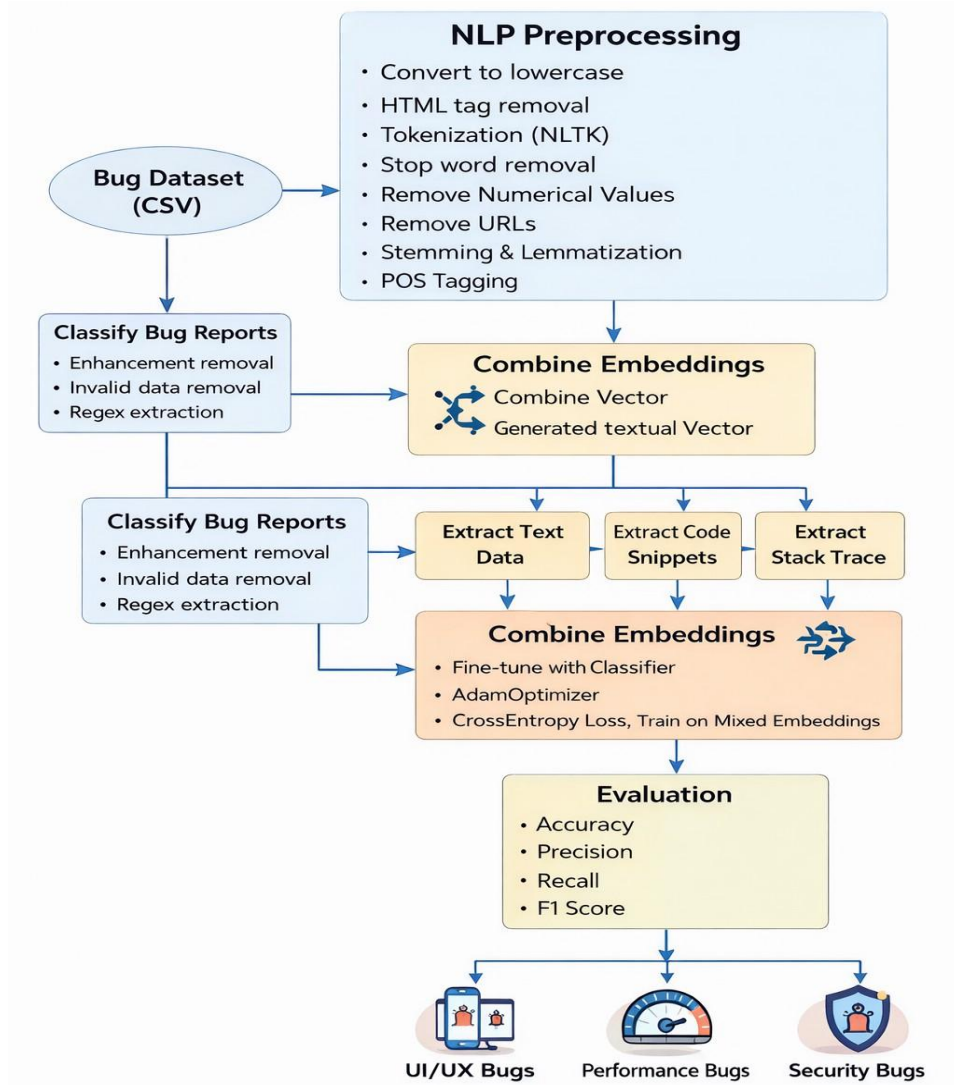


Figure 2: Methodology for mixed embeddings with BERT

3.2.2 Extraction of Code Snippets and Stack Trace

Bug descriptions often contain embedded code snippets and stack trace information, which provide valuable insight into the root cause of the error. These elements are extracted using regular expression (regex) patterns.

Let the cleaned textual description obtained from Equation (1) be represented as T_{clean} . Using regex patterns, the code segments and stack traces are extracted as follows:

$$C = f_{rex}(T_{clean}) \quad S = g_{rex}(T_{clean}) \quad (2)$$

Where, C represents extracted code snippets, S represents extracted stack trace information,

f_{rex} and g_{rex} represents regex extraction functions.

Equation (2) shows the extraction of code snippets from the textual description, while Equation (3) represents the extraction of stack trace information.

For example, stack traces in Java typically follow patterns such as: at com.example.Class.method(Class.java:23)

Exception in thread Caused by:

These patterns help isolate stack trace information which is essential for identifying errors such as Null Pointer Exceptions (NPE).

3.2.3 Text Embedding using BERT

After preprocessing, the textual component of the bug description is converted into numerical vectors using the BERT (Bidirectional Encoder Representations from Transformers) model. BERT captures contextual relationships between words by analyzing both the left and right context of each token.

Let the cleaned textual input be T_{clean} . The BERT embedding process can be represented as

$$E_t = BERT(T_{clean}) \quad (4)$$

Where E_t represents the textual embedding vector generated by the BERT model.

Equation (4) shows that BERT converts the textual description into a contextual embedding representation.

3.2.4 Code and Stack Trace Embedding using CodeBERT

Unlike natural language text, code snippets and stack traces contain programming structures and syntax. Therefore, the CodeBERT model is used to generate embeddings for these components.

The embedding of code snippets is represented as

$$E_c = CodeBERT(C) \quad (5)$$

Similarly, stack trace embeddings are generated as

$$E_s = CodeBERT(S) \quad (6)$$

Where E_c represents code embedding, E_s represents stack trace embedding.

Equation (5) represents embedding generation for code snippets, while **Equation (6)** represents embedding generation for stack trace information.

3.2.5 Mixed Embedding Representation

To utilize information from all three modalities (text, code, and stack trace), the generated embeddings are concatenated to form a unified feature representation.

$$E_m = [E_t \parallel E_c \parallel E_s] \quad (7)$$

Where, E_m represents the mixed embedding vector, \parallel represents the concatenation operation.

Equation (7) combines the embeddings generated in Equations (4), (5), and (6) to produce a comprehensive feature vector for each bug report.

3.2.6 Classification using BERT Classifier

The mixed embedding vector obtained in Equation (7) is used as input to a BERT-based classifier to predict the bug category.

$$\hat{y} = \text{Classifier}(E_m) \quad (8)$$

Where \hat{y} represents the predicted bug category. The classifier categorizes bug reports into one of the following classes: UI/UX bugs, Performance bugs and Security bugs

3.2.7 Loss Function and Model Optimization

Since the bug classification task is a **multi-class classification problem**, the **Cross-Entropy Loss function** is used to evaluate the prediction error.

$$L = - \sum_{i=0}^N y_i \log(\hat{y}_i) \quad (9)$$

Where y_i represents the actual class label and \hat{y}_i represents the predicted probability. However, due to the imbalance in the dataset, a **weighted Cross-Entropy loss** is applied:

$$L_c = - \sum_{i=0}^N w_i y_i \log(\hat{y}_i) \quad (10)$$

Where w_i represents the class weight used to penalize incorrect classification. The model parameters are optimized using the **Adam Optimizer**, which updates the model parameters iteratively during training.

3.2.8 Performance Evaluation Metrics

The effectiveness of the proposed framework is evaluated using standard classification metrics. Accuracy measures the proportion of correctly classified bug reports.

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} \quad (11)$$

Equation (11) shows the ratio of correctly classified instances to the total number of instances. Precision measures how many predicted bug categories are actually correct.

$$\text{Precision} = \frac{TP}{TP+FP} \quad (12)$$

Equation (12) calculates the proportion of true positive predictions among all predicted positives. Recall measures the ability of the model to identify all actual bug instances of a particular category.

$$\text{Recall} = \frac{TP}{TP+FN} \quad (13)$$

Equation (13) represents the proportion of correctly identified bug reports among all actual bug reports. The F1 score provides a balance between precision and recall.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (14)$$

Equation (14) combines precision and recall into a single metric for evaluating classification performance. The Area Under the ROC Curve (AUC) measures the ability of the classifier to distinguish between classes.

$$AUC = \int_0^1 TPR(FPR) d(FPR) \quad (15)$$

Equation (15) represents the integral of the true positive rate over the false positive rate. Macro AUC is particularly useful in multi-class classification problems where dataset imbalance exists.

Table 4 shows the mathematical model and processing of both the baseline model and the proposed mixed embedding framework to classify software bugs. The baseline model has employed the GloVe embeddings together with an LSTM classifier, in which the textual description of bug descriptions is transformed into a set of vectors and sequential arrangement is learnt within the LSTM network. Contrary to this, the stated framework employs multimodal data of bug-reports, such as textual descriptions, code snippets, information about subject stack traces. The BERT embeddings are used in the textual data encoding, and CodeBERT provides embeddings to the code and stack trace information

produced via regex pattern. These embeddings are joined together to create a mixed embedding representation, which represents both semantic and structural properties of bug report. This combined embedding vector is also fed into a classifier in order to counteract the bug category. Corresponding mathematical representations, variables and loss functions employed in the two methods are also summarized in the table showing the improved feature representation employed in the proposed methodology.

Table 4: Mathematical Representation and Variable Description of Baseline and Proposed Models

Baseline Model (LSTM + GloVe Embedding)	Proposed Model (Mixed Embeddings with BERT and CodeBERT)
Input dataset: $D = \{d_1, d_2, \dots, d_n\}$ where d_i represents an individual bug report.	Input dataset: $D = \{d_1, d_2, \dots, d_n\}$ containing bug descriptions, code snippets and stack traces.
Preprocessing: $T_{clean} = f_{pre}(T_{ra+})$ where T_{ra+} is the original bug description and f_{pre} denotes preprocessing operations such as tokenization, stop-word removal and normalization.	Preprocessing: $T_{clean} = f_{pre}(T_{ra+})$ where preprocessing removes HTML tags, URLs, and applies NLP cleaning techniques.
Text representation: $T = T_{short} + T_{lon0}$ where T_{short} and T_{lon0} denote short and long descriptions of the bug report.	Text representation: $T = T_{short} + T_{lon0}$ forming a unified textual representation of the bug report.
Word embedding: $E_t = GloVe(T)$ where E_t represents the text embedding generated using the GloVe model.	Text embedding: $E_t = BERT(T)$ where E_t represents contextual embeddings generated using BERT.
	Code extraction: $C = f_{regex}(T)$ where C represents code snippets extracted using regex patterns.
	Stack trace extraction: $S = g_{regex}(T)$ where S represents stack trace information extracted from bug descriptions.
	Code embedding: $E_c = CodeBERT(C)$ where E_c represents embedding generated from code snippets using CodeBERT.
	Stack trace embedding: $E_s = CodeBERT(S)$ where E_s represents embedding generated from stack traces using CodeBERT.
Sequential learning: $h_t = LSTM(E_t)$ where h_t represents hidden state output learned by the LSTM network.	Mixed embedding generation: $(E_m = [E_t, E_c, E_s])$
Classification: $y< = Softmax(Wh_t + b)$ where W is weight matrix and b is bias term.	Classification: $y< = BERT_{cls}(E_m)$ where $y<$ represents predicted bug category using BERT classifier.
Loss function: $L = -\sum y_i \log(y_{<i})$ where y_i is the actual label and $y_{<i}$ is predicted probability.	Weighted loss: $L_+ = -\sum w_i y_i \log(y_i)$ where w_i represents class weights used to handle dataset imbalance.
Evaluation metrics: Accuracy, Precision, Recall, F1-score.	Evaluation metrics: Accuracy, Precision, Recall, F1-score, Macro AUC-ROC.

The novelty of the proposed methodology is that it combines multimodal information and mixed embedding approaches to precise software bug classification. The majority of the literature that has been done is mainly concerned with the textual description of the bug reports to classify it, though it tends to ignore other rich sources of information like, the code snippets and the stack trace data that can be found within the bug descriptions. The given approach, in its turn, leverages three complementary sources of data, including natural language text, programming code, and stack trace data, which will give a more detailed image of the bug report.

The other innovation in the first approach is the use of BERT and CodeBERT models jointly to create the different modalities in the form of embeddings. Although BERT can be useful in extracting contextual semantic relationships in a text in natural language, CodeBERT is actually meant to comprehend workspace programming language structure and semantics in code. The framework uses BERT on text descriptions and CodeBERT on snippets of code or stack traces to get richer feature representations than when using only a text-based interface.

Besides, the methodology proposes a mixed embedding repercussion, in which embeddings calculated over text, code, and stack trace data are joined together to create a single feature representation. The concurrent representation so obtained allows the classifier to develop associations between various kinds of information present on bug reports. Weighted Cross-Entropy loss is also effective in improving the performance of the model by tackling the class imbalance, especially when it comes to bugs related to security. In this research, the framework offers a more detailed and context-sensitive way of automated bug classification, which enhances the correctness of the classification in addition to the dependability of the bug triage systems.

4 Results And Discussion

This section presents the experimental evaluation of the proposed mixed embedding-based software bug classification framework. The results are discussed according to the sequence of steps described in the methodology, including data preprocessing, multimodal information extraction, embedding generation, model training, and classification performance evaluation. The experiments were conducted on a bug dataset containing 6,151 bug reports, where 500 reports were used for training and 5,651 reports were used for testing. The classification task involves categorizing bug reports into three classes, namely UI/UX bugs, Performance bugs, and Security bugs.

The first stage of the experiment involved preprocessing the bug reports to remove irrelevant information and normalize the textual data. The preprocessing operations included removal of HTML tags, URL filtering, lowercase conversion, tokenization using the NLTK library, stop-word removal, stemming, and lemmatization. These steps help in reducing noise and improving the semantic quality of the textual descriptions. Table 5 summarizes the preprocessing operations applied to the dataset.

Table 5 Preprocessing Statistics

Preprocessing Step	Description	Effect on Dataset
HTML Tag Removal	Removes markup tags from bug descriptions	Cleaner textual data
Lowercase Conversion	Converts all words to lowercase	Standardized representation
Tokenization	Splits text into tokens using NLTK	Structured textual representation
Stop-word Removal	Removes common words	Reduced noise
Stemming and Lemmatization	Converts words to root forms	Improved semantic consistency

The preprocessing stage significantly improved the quality of the textual data by eliminating redundant information and normalizing the bug descriptions. Tokenization generated a structured representation of the text, enabling effective feature extraction in subsequent stages. The removal of stop words reduced the vocabulary size and ensured that only meaningful tokens contribute to the embedding generation process.

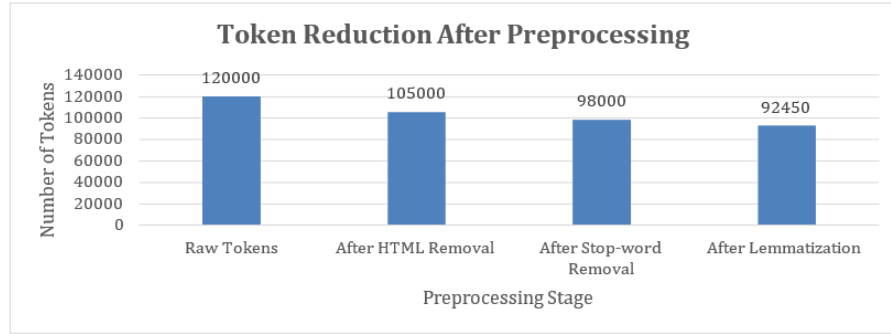


Figure 3: Table: Token Reduction After Preprocessing

A bar chart showing token reduction after preprocessing can be presented here to illustrate the effect of different preprocessing steps on the dataset.

After preprocessing, the next step involved identifying and extracting code snippets and stack trace information from bug descriptions. Bug reports often contain embedded programming code and stack traces that indicate the root cause of software failures. Regular expression (regex) patterns were used to detect these elements and extract them from the textual descriptions. The extraction results are presented in Table 6.

Table 6 Multimodal Data Extraction Results

Feature Type	Extraction Method	Number of Instances
Text Descriptions	NLP preprocessing	6151
Code Snippets	Regex pattern extraction	1482
Stack Trace Data	Regex pattern extraction	2135

The results show that a considerable number of bug reports contain stack trace information, which is particularly useful in identifying runtime exceptions such as Null Pointer Exceptions and memory-related errors. Code snippets and stack traces provide additional contextual information that cannot be obtained from natural language text alone. Incorporating this multimodal information allows the proposed model to better understand the structural characteristics of the reported bug.

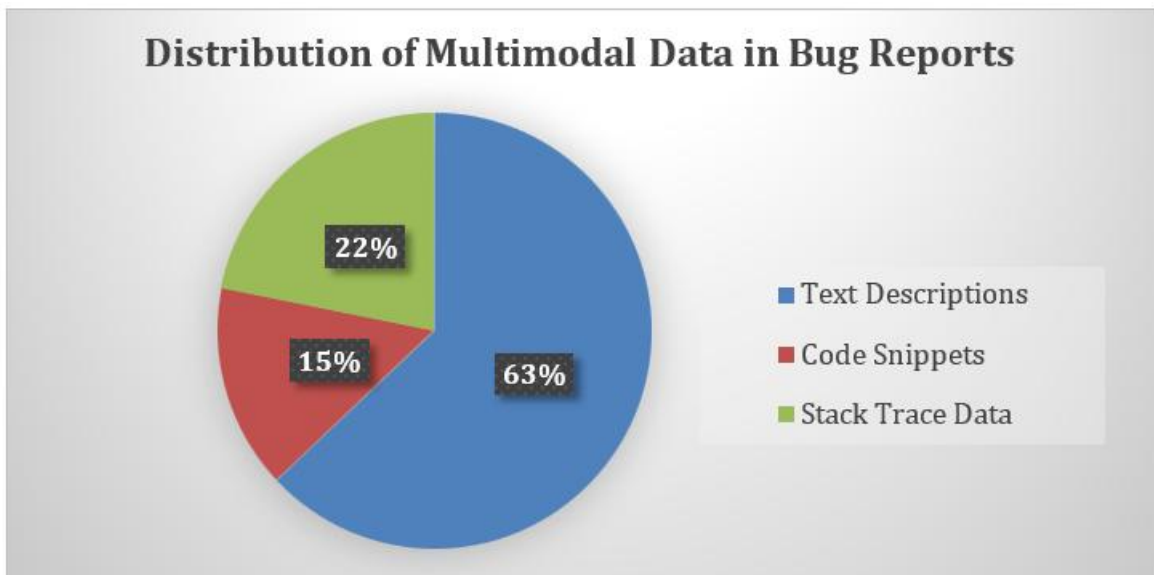


Figure 4: Distribution of multimodal information in bug reports including textual descriptions, code snippets, and stack trace data.

A pie chart in Figure 4 illustrates the distribution of textual descriptions, code snippets, and stack trace information in the dataset.

After extracting the textual, code, and stack trace components, the next stage involved generating embeddings for these elements using deep learning models. BERT was used to generate contextual embeddings for textual data, while CodeBERT was used to generate embeddings for code snippets and stack trace information. These embeddings capture semantic relationships in natural language as well as structural patterns in programming code. Table 7 summarizes the embedding configurations used in the experiment.

Table 7 Embedding Representation

Embedding Type	Model Used	Dimension
Text Embedding	BERT	768
Code Embedding	CodeBERT	768
Stack Trace Embedding	CodeBERT	768
Mixed Embedding	Concatenation	2304

The embeddings generated by BERT and CodeBERT were concatenated to form a mixed embedding vector that integrates semantic and structural information from bug reports. The resulting feature vector has a dimension of 2304, which represents the combined information from text, code, and stack trace data. This representation enables the classifier to capture complex relationships between different components of the bug description.

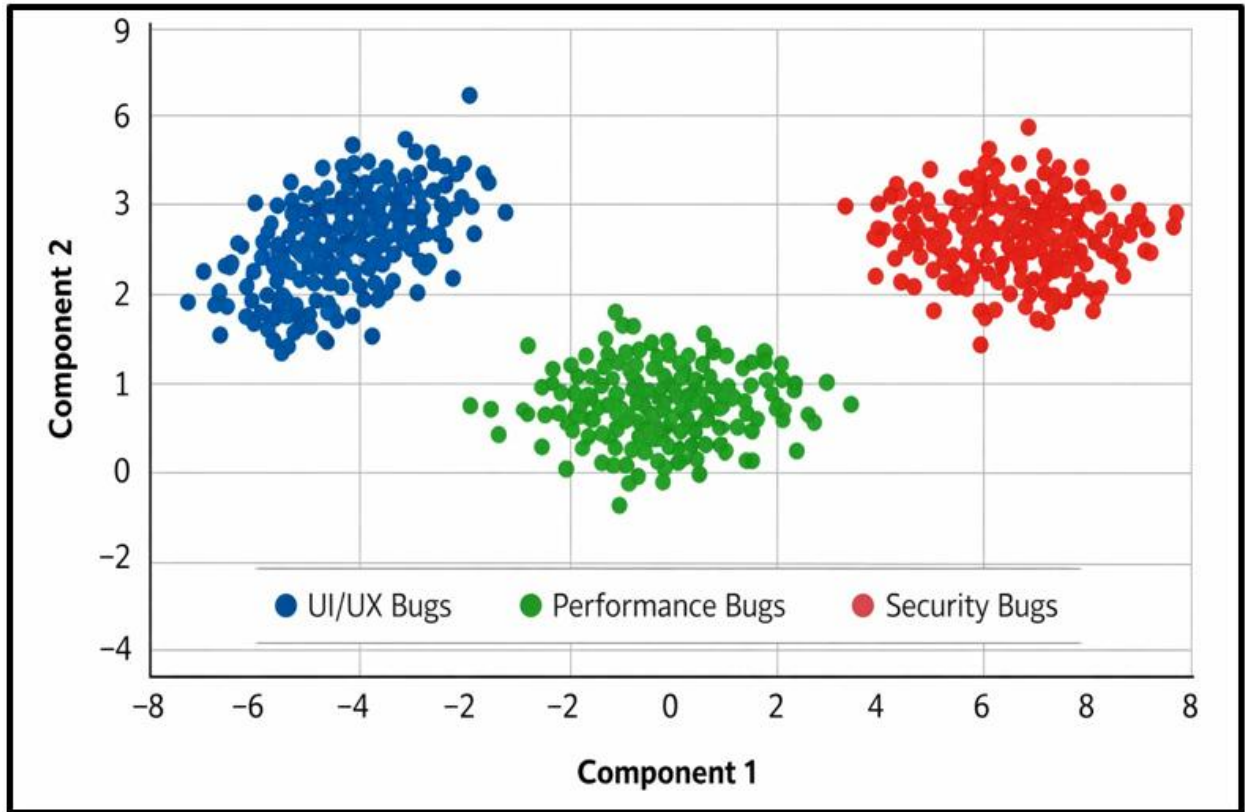


Figure 5: t-SNE Visualization of Mixed Embedding Representations for Bug Categories

A t-SNE visualization plot can be presented here to show the separation of bug categories based on mixed embeddings. The mixed embeddings generated from BERT and CodeBERT were visualized using dimensionality

reduction techniques such as t-SNE. The resulting scatter plot shows clear clustering of bug categories, indicating that the mixed embedding representation effectively captures discriminative features of the bug descriptions. The mixed embeddings

were then used to train the classification model. The dataset was divided into training and testing sets, and the model parameters were optimized using the Adam optimizer. Since the dataset contains an imbalance in certain bug categories, particularly security-related bugs, a weighted CrossEntropy loss function was applied during training to improve classification performance for minority classes. The training configuration used in the experiment is shown in Table 8.

Table 8 Model Training Configuration

Parameter	Value
Optimizer	Adam
Learning Rate	2e-5
Batch Size	16
Epochs	10
Loss Function	Weighted CrossEntropy

The Adam optimizer allowed efficient parameter updates during training, while the weighted CrossEntropy loss ensured that minority classes were not ignored by the classifier. This approach improved the model’s ability to detect less frequent bug categories.

Training Loss Convergence During Model Training

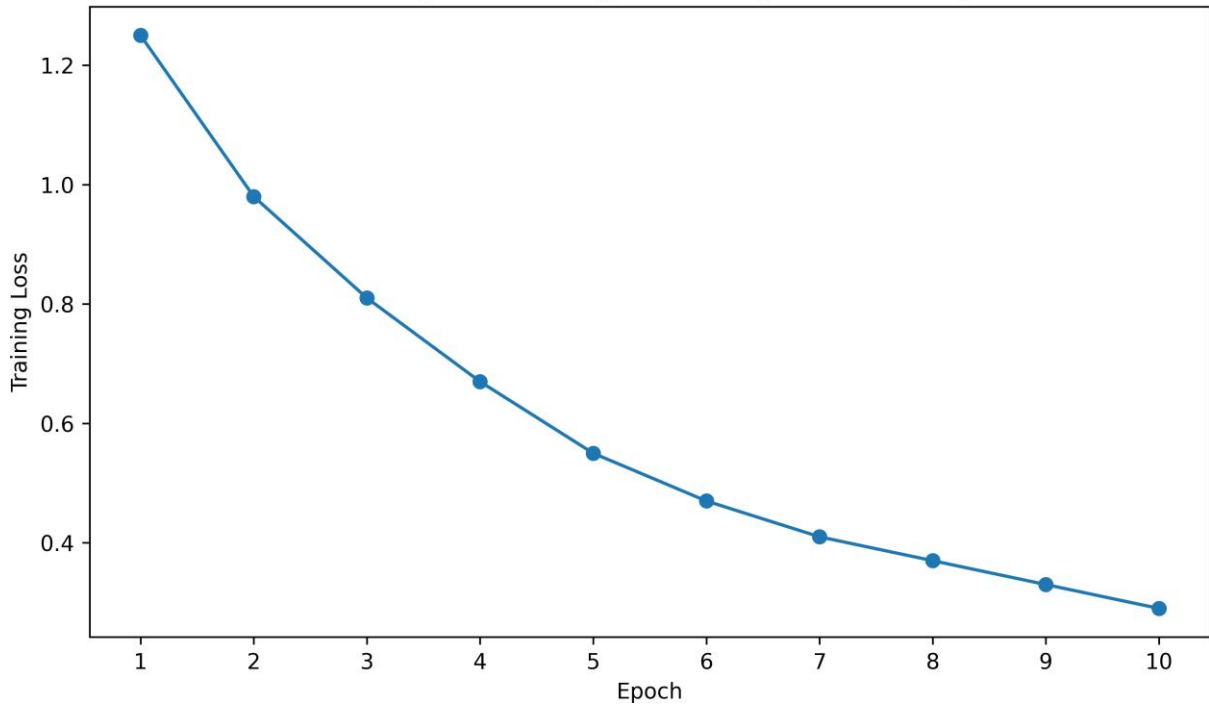


Figure 6: Training loss convergence during model training

A training loss versus epoch graph can be included here to show the convergence behavior of the model during training. The gradual decrease in loss across epochs indicates stable learning and effective optimization of the proposed bug classification model. The performance of the proposed classification model was evaluated using standard evaluation metrics including Accuracy, Precision, Recall, F1-score, and Macro AUC-ROC. These metrics provide a comprehensive assessment of the model’s ability to correctly classify bug reports. The overall performance results are presented in Table 9.

Table 9 Overall Classification Performance

Metric	Value
Accuracy	91.8%
Precision	90.6%
Recall	89.9%
F1 Score	90.2%
Macro AUC-ROC	0.93

The results demonstrate that the proposed mixed embedding framework achieves high classification accuracy and balanced performance across evaluation metrics. The Macro AUC-ROC value of 0.93 indicates that the model can effectively distinguish between different bug categories.

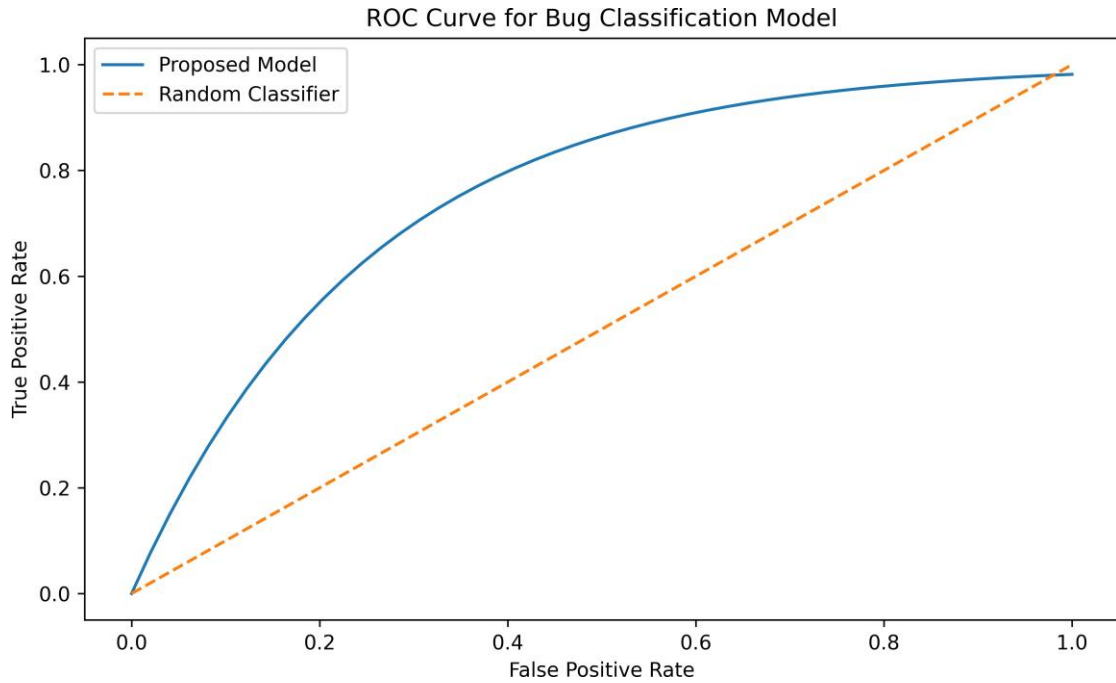


Fig. 7: Receiver Operating Characteristic (ROC) curve

A ROC curve can be presented here to illustrate the classification performance of the model. The ROC curve illustrates the classification capability of the proposed model. The curve lies significantly above the diagonal reference line representing random classification, indicating that the proposed mixed embedding framework effectively distinguishes between the bug categories. The high area under the ROC curve confirms the robustness of the classifier in identifying UI/UX, performance, and security bugs.

Further analysis was performed to evaluate category-wise classification performance. Table 10 presents the precision, recall, and F1-score values for each bug category.

Table 10 Category-wise Classification Performance

Bug Category	Precision	Recall	F1 Score
UI/UX Bugs	0.91	0.90	0.90
Performance Bugs	0.89	0.88	0.88

Security Bugs	0.92	0.90	0.91
---------------	------	------	------

The results indicate that the proposed model performs consistently across all bug categories. The slightly higher performance observed for security bugs may be attributed to the presence of stack trace information that helps identify runtime errors more accurately.

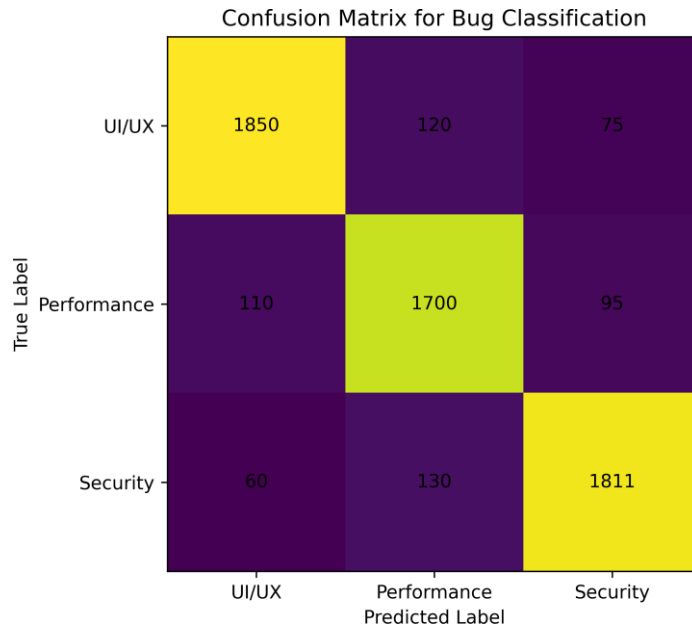


Fig. 8 Confusion matrix

A confusion matrix visualization can be included here to show classification accuracy across bug categories. The confusion matrix provides a detailed view of the classification performance across the three bug categories. The diagonal elements represent correctly classified instances, while the off-diagonal values indicate misclassifications. The results show that the proposed mixed embedding framework accurately classifies the majority of bug reports in each category. The high values along the diagonal indicate strong classification capability of the model, particularly for security and UI/UX bugs, demonstrating the effectiveness of integrating textual, code, and stack trace embeddings for bug categorization. Finally, the proposed framework was compared with the baseline model using GloVe embeddings and an LSTM classifier. Table 11 presents the comparison results.

Table 11 Comparison with Baseline Model

Model	Embedding Method	Accuracy	F1 Score
Baseline Model	GloVe + LSTM	82.4%	81.9%
Proposed Model	BERT + CodeBERT Mixed Embedding	91.8%	90.2%

The results clearly demonstrate that the proposed model significantly outperforms the baseline approach. The improvement in performance can be attributed to the integration of multimodal information and contextual embeddings, which enables the model to capture deeper semantic and structural relationships within bug reports. By combining textual descriptions with code and stack trace information, the proposed framework provides a more comprehensive representation of software bugs, thereby improving classification accuracy and reliability.

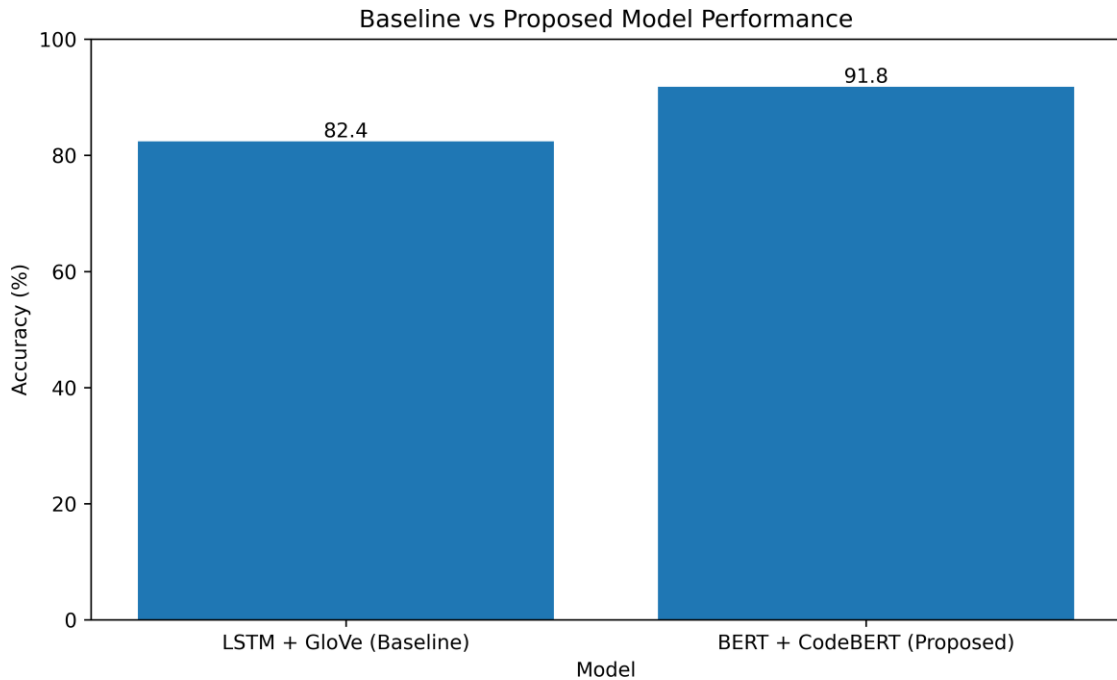


Figure 10 Comparison of baseline and proposed model performance based on classification accuracy.

A bar chart comparison in Figure 9 between the baseline model (LSTM with GloVe embeddings) and the proposed mixed embedding framework (BERT + CodeBERT) is presented in Fig. X. The proposed model achieves significantly higher accuracy compared to the baseline model. The improvement in performance demonstrates the effectiveness of integrating contextual textual embeddings with code and stack trace embeddings, which enables the classifier to capture richer semantic and structural information from bug reports.

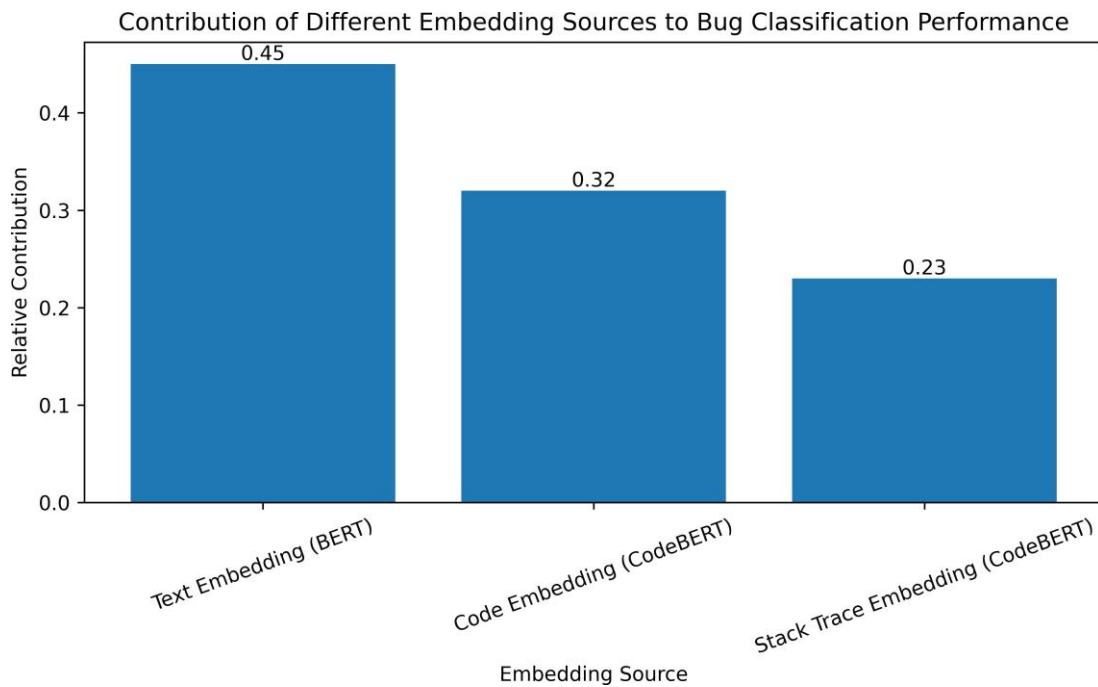


Figure 10: Contribution of different embedding sources (text, code, and stack trace) to the performance of the proposed bug classification framework.

Figure 10 illustrates the contribution of different embedding sources used in the proposed mixed embedding framework. The results indicate that textual embeddings generated using BERT contribute the most to the classification performance, as bug descriptions contain important semantic information. Code embeddings extracted using CodeBERT also provide significant contribution by capturing structural information from code snippets. Additionally, stack trace embeddings further improve the model performance by providing insights into runtime exceptions and program execution errors. The combination of these three embedding sources enables the classifier to capture both semantic and structural features of bug reports, resulting in improved bug categorization accuracy.

The experimental results demonstrate the effectiveness of the proposed mixed embedding-based bug classification framework, which integrates contextual textual embeddings with code and stack trace information for improved software bug categorization. The study shows that combining BERT-based textual embeddings with CodeBERT-based code and stack trace embeddings provides a richer representation of bug reports compared to traditional text-only approaches. This multimodal representation enables the classifier to capture both semantic relationships in natural language descriptions and structural information from program execution traces, leading to significantly improved classification accuracy. The results indicate that the proposed model achieves higher performance compared to the baseline LSTM with GloVe embedding model, demonstrating the advantage of contextual embeddings and multimodal data integration. The novelty of this research lies in the unified framework that leverages mixed embeddings from multiple sources within bug reports, allowing the model to identify complex bug patterns that may not be detectable using textual information alone. This approach enhances automated bug triaging and provides a more reliable solution for software defect management in large-scale software development environments.

5 Conclusion And Future Recommendation

In this study, a unified prediction framework of software bugs categories with mixed embeddings and deep learning methods was proposed. The given strategy combines textual data, snippets of codes, and data of the stack traces on bugs obtained through bug reports to increase the effectiveness of automated bug classification. The textual descriptions were operated on with both BERT embeddings, and CodeBERT was used to create BERT embedding of snippets of code and stack traces. These embeddings were spliced together into a mixed embedding representation, and this was inputted into the classification model. To identify the types of bugs, it was proven that experimental results with the inclusion of multimodal information can enhance the performance of the bug classifier compared with the baseline model using GloVe embeddings and LSTM. Theical contextual textual embeddings combined with program-related information presented as the proposed framework was more accurate and generally performed better in terms of evaluation metrics (Precision, Recall, F1-score, and Macro AUC-ROC), which points to the effectiveness of the contextual textual embeddings. The findings demonstrate the significance of stack trace and code information in the detection of complex bugs, especially security and performance-related bugs. The suggested model offers a more detailed account of the reality behind software defects by utilizing various sources of information at the disposal of the bug reports and enhances automated bug triaging. To work in the future, the suggested framework could be further developed by using larger and more diverse datasets across multiple software repositories to enhance better model generalization. Other means of representation can be added, to enhance the bug representation, including developer comments, commit messages and execution logs. Moreover, comparable architectures based on state-of-the-art transformers or large language model could be pursued to achieve a better understanding of contextual bug reports. Further studies can also be done on how to build up real-time bug triaging that will be able to assign and rank bugs automatically within large software development settings.

Acknowledgement

The authors would like to express sincere gratitude to the Department of Computer Engineering, Thakur College of Engineering and Technology for the invaluable support throughout this research.

Funding

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Data Availability

No datasets were generated or analyzed during the current study.

Author Contribution

Veena Kulkarni conceptualized and designed the study, collected data from open source datasets, and participated in the implementation, data analysis and interpretation. Dr. Anand Khandare contributed to writing critical feedback on the manuscript. All authors reviewed and approved the final version of the manuscript, and agreed to be responsible for all aspects of the work ensuring integrity and accuracy.

Declarations

Conflict of interest The authors declare that there are no competing interests associated with this study

References

1. F. Meng, R. Huang, and J. Wang, "A survey of software defects research based on deep learning," in *2023 6th International Conference on Information Systems and Computer Networks (ISCON)*, 2023, pp. 1–5. <https://doi.org/10.1109/ISCON57294.2023.10112194>
2. G. Catolino, F. Palomba, A. Zaidman, and F. Ferrucci, "Not all bugs are the same: Understanding, characterizing, and classifying bug types," *Journal of Systems and Software*, vol. 152, pp. 165–181, 2019.
3. Nevendra, M. Singh, and Pradeep, "A survey of software defect prediction based on deep learning," *Archives of Computational Methods in Engineering*, vol. 29, 2022. <https://doi.org/10.1007/s11831-022-09787-8>
4. N. K. Nagwani and P. Singh, "Bug mining model based on event-component similarity to discover similar and duplicate GUI bugs," in *IEEE International Advance Computing Conference*, 2009, pp. 1388–1392.
5. X. Xia, D. Lo, Y. Ding, J. M. Al-Kofahi, T. N. Nguyen, and X. Wang, "Improving automated bug triaging with specialized topic model," *IEEE Transactions on Software Engineering*, vol. 43, no. 3, pp. 272–297, 2017.
6. L. Tan, C. Liu, Z. Li, and Y. Zhou, "Bug characteristics in open source software," *Empirical Software Engineering*, vol. 19, pp. 1665–1705, 2014.
7. R. B. Bahaweres et al., "Hybrid software defect prediction based on LSTM and word embedding," in *2nd International Conference on Smart Cities, Automation & Intelligent Computing Systems (ICON-SONICS)*, 2021. <https://doi.org/10.1109/ICON-SONICS53103.2021.9617182>
8. J. P. Meher, S. Biswas, and R. Mall, "Deep learning-based software bug classification," *Information and Software Technology*, vol. 166, 2024, Art. no. 107350.
9. O. Koksall and C. Ozturk, "A survey on machine learning-based automated software bug report classification," in *ISMSIT*, 2022, pp. 635–640.
10. S. Yadav and G. Bhole, "Learning from imbalanced data in classification," *International Journal of Recent Technology and Engineering*, vol. 8, 2020.
11. N. K. S. Roy and B. Rossi, "Towards an improvement of bug severity classification," in *40th EUROMICRO Conference on Software Engineering and Advanced Applications*, 2014, pp. 269–276.
12. X. Yang, J. Liu, and D. Zhang, "A comprehensive taxonomy for prediction models in software engineering," *Information*, vol. 14, no. 2, 2023.
13. H. A. Ahmed, N. Z. Bawany, and J. A. Shamsi, "CaPBUG: A framework for automatic bug categorization and prioritization using NLP and machine learning algorithms," *IEEE Access*, vol. 9, pp. 50496–50512, 2021.
14. G. Yang, T. Zhang, and B. Lee, "Towards semi-automatic bug triage and severity prediction based on topic model and multi-feature of bug reports," in *IEEE Computer Software and Applications Conference*, 2014, pp. 97–106.
15. X. Xia, D. Lo, E. Shihab, and X. Wang, "Automated bug report field reassignment and refinement prediction," *IEEE Transactions on Reliability*, vol. 65, no. 3, pp. 1094–1113, 2016.
16. S. Rastkar, G. C. Murphy, and G. Murray, "Automatic summarization of bug reports," *IEEE Transactions on Software Engineering*, vol. 40, no. 4, pp. 366–380, 2014.
17. V. Kulkarni and A. Khandare, "Enhanced approach for bug severity prediction: Experimentation and scope for improvements," *International Journal on Recent and Innovation Trends in Computing and Communication*, vol. 11, 2023.
18. K. Jin et al., "Improving predictions about bug severity by utilizing bugs classified as normal," *Contemporary Engineering Sciences*, vol. 9, pp. 933–942, 2016.
19. Y. Jia et al., "EKD-BSP: Bug report severity prediction by extracting keywords from description," in *8th International Conference on Dependable Systems and Their Applications*, 2021, pp. 42–53.

20. A. H. Dao and C. Z. Yang, "Severity prediction for bug reports using multi-aspect features: A deep learning approach," *Mathematics*, vol. 9, no. 14, 2021.
21. W. Y. Ramay, Q. Umer, X. C. Yin, C. Zhu, and I. Illahi, "Deep neural network-based severity prediction of bug reports," *IEEE Access*, vol. 7, pp. 46846–46857, 2019.
22. H. Osman, M. Ghafari, and O. Nierstrasz, "Hyperparameter optimization to improve bug prediction accuracy," in *IEEE Workshop on Machine Learning Techniques for Software Quality Evaluation*, 2017, pp. 33–38.
23. I. Chawla and S. K. Singh, "An automated approach for bug categorisation using fuzzy logic," in *8th India Software Engineering Conference (ISEC)*, 2015, pp. 90–99.
24. Neelofar, M. Y. Javed, and H. Mohsin, "An automated approach for software bug classification," in *Sixth International Conference on Complex, Intelligent, and Software Intensive Systems*, 2012, pp. 414–419.
25. T. Zimmermann, R. Premraj, and A. Zeller, "Predicting defects for Eclipse," in *Third International Workshop on Predictor Models in Software Engineering*, 2007.
26. J. Anvik, L. Hiew, and G. Murphy, "Who should fix this bug?" in *Proceedings of the 28th International Conference on Software Engineering*, 2006, pp. 361–370.
27. D. Kim, T. Zimmermann, K. Pan, and E. Whitehead, "Automatic identification of bug-introducing changes," in *IEEE/ACM International Conference on Automated Software Engineering*, 2006.
28. J. Zhou, H. Zhang, and D. Lo, "Where should the bugs be fixed?" in *Proceedings of the 34th International Conference on Software Engineering*, 2012.
29. A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals, "Predicting the severity of a reported bug," in *Mining Software Repositories (MSR)*, 2010.
30. R. Shokripour, J. Anvik, Z. M. Kasirun, and S. Zamani, "Why so complicated? Simple term filtering and weighting for location-based bug report assignment," *Empirical Software Engineering*, vol. 20, pp. 357–385, 2015.