



# Next-Generation AI-Assisted Bug Tracking and Automated Code Analysis System

Priya Pandey<sup>1</sup>, \*Aniket Kumar<sup>2</sup>, Aditya Singh<sup>3</sup> and Shruti Kumari<sup>4</sup>

<sup>1,2,3,4</sup>Department of Computer Science and Engineering, Galgotias University, Greater Noida, Uttar Pradesh, India

<sup>1</sup>priyapandeyrmp@gmail.com

<sup>2</sup>aniketkumar1623@gmail.com

<sup>3</sup>vikku6285@gmail.com

<sup>4</sup>shrutir28@gmail.com

**Abstract.** Bug tracking is a time-consuming activity in current software development, which generally involves manual bug reporting, bug triage, and bug propagation. In large and fast-evolving projects, this manual dependency causes slowing down bug resolution, duplicating reports, inconsistency in labeling severity, and overburdening developers. This paper introduces an AI-assisted bug tracking and automated code analysis system that aims at lowering human participation in the process of defect management. The system combines widely used static code analysis applications – ESLint and SonarQube – with our custom-built parsing modules, an XGBoost-based severity classifier, a hybrid Sentence-BERT duplicate detector, and a retrieval augmentation-based fix recommendation module, to automatically predict the severity of bugs, approximate their priority, identify duplicate bug reports, and provide fix recommendations from past fixed cases. The experimental results show that the proposed system produces 91.4% classification accuracy of bugs, 65% reduction of triage time (42.3 minutes to 14.8 minutes), and 23% reduction of duplicate report submissions which surpasses two traditional systems, Jira, Bugzilla, and YouTrack, and machine learning baselines such as Naive Bayes (78.2%), SVM (82.6%), CNN (85.1%), and Distil BERT (88.7), and all results are statistically significant  $p < 0.001$ . The main contribution of this work is the development of a fully-integrated end-to-end pipeline that unifies the automated static code analysis and the AI-supported triage done by the human developer into a single CI-CD deployable system that provides measurable success in terms of the quality of software defect resolution and developer productivity compared to the traditional manual tracking approach.

**Keywords:** Bug Tracking System, Static Code Analysis, ESL int, SonarQube, Artificial Intelligence

## 1 INTRODUCTION

Software systems have grown in complexity at a remarkable pace, and distributed development teams and continuous integration pipelines are associated with higher defect rates [1][2]. The traditional bug tracking software, such as Jira, Bugzilla, and YouTrack, provides structured workflow management of defects but relies heavily on manual processes of detecting, classifying, and prioritizing bugs. A variety of

underlying inefficiencies are inherent in such a manual dependency: severity ratings are not homogeneous, reporting is redundant, slow response, and a larger workload among the developers [3][4]. The large-scale projects are known to worsen these inefficiencies because the bugs can be reported by hundreds of people per day. Automatic code scanners like ESLint and SonarQube will provide reports containing the potential faults during the early phase of the development process, but the data will not be connected to the bug tracking process and must be manually interpreted and integrated into the bug tracking workflow. The advancements in machine learning have enabled the provision of prospective solutions to the solitary bug management operations: severity prediction [2], duplicate identification [5], and automated triage [3]. However, there are three severe flaws in the existing solutions: they are not end-to-end automated, they are not integrated into production development environments, and they are not constantly learning to respond to mutating codebases.

In this paper, the limitations will be addressed by proposing the concept of an integrated AI-based bug tracker and automatic code analysis as a single platform, which will rely on the use of the tools of static analysis, the developed parsers, and machine-learning models. We can automatically detect bugs, standardize analysis reports between tools, predict severity and priority through trained ML models, give duplicate recommendations based on semantic similarity, and recommend fixes by providing context. Although earlier systems do individual subtasks one by one - DeepTriage [7] predicts severity through ensemble methods, but does not actually integrate a static analysis; BugDedup [5] uses Sentence-BERT to do duplicate detection, but again does not actually combine with a bug life cycle management system; the current tools such as SonarQube and ESLint generate reports that are completely disconnected with a bug lifecycle management system. As shown in Table 1, to the best of our knowledge, our system is the first to combine ingestion of static analysis, severity classification based on ML, and hybrid semantic duplicate detection with retrieval augmented fix recommendation into a single deployable production pipeline [6] [7].

The main contributions of the work are organized as follows:

**Algorithms Contributions:** –

- A hybrid similarity fusion model of duplicate detection, which is the main contribution to the algorithm, that combines a Sentence-BERT semantic embedding (weight  $a=0.7$ ), the Jaccard file overlap ( $b=0.2$ ), and component matching ( $g=0.1$ ) with learned ratios.
- A multi-class severity classifier using XGBoost with 91.4% accuracy, which was better than Naive Bayes, SVM, CNN, and DistilBERT baselines using the same dataset.
- A fix recommendation module based on retrieval with augmentation with templates matching and SBERT-indexed historical fix retrieval. – Evaluation on 95,400 bug reports on four large open-source projects statistically confirmed to have made substantial improvements over the manual baselines.

**Architectural Contributions:** –

- An end-to-end pipeline that is production-ready and that combines ESLint, SonarQube, and custom parsers together with AI-based triage services and creates one single bug tracking environment with complete CI/CD/ supported features.

- A microservice architecture that is modular and allows the update of the components independently, and provides a continuous learning process without bringing the system to a halt.

## 2 RELATED WORK

Bug tracking and defect management have been widely studied in four areas: traditional bug tracking systems [3][15][16], static code analysis [3][13], ML-based severity and priority prediction [6][7][8], and duplicate bug detection [4][5][14]. Traditional systems like Bugzilla, Jira, and YouTrack offer structured management of the workflow but are purely based on manual triage, leading to inconsistent severity labels and redundant reports [2][4]. Static analysis tools – ESLint, SonarQube, and PMD [3][13] – identify code smells with vulnerabilities early, but are completely disconnected from the bug tracking lifecycle with no automated triage functionality [5]. ML-based severity prediction models, including Naive Bayes [6], ensemble approaches [7][8], and a ranker with CNN [2], deal well with individual classification problems but work on static offline datasets rather than being part of continuous delivery and retraining.

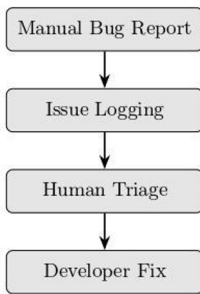
Duplicate detection approaches – TF-IDF and LSI [3][4] for lexical overlap, Sentence-BERT [5] for semantic similarity – increase the accuracy of detection but lack structural signals such as file overlap and component matching that work highly in addition to semantic evidence. As summarized in Table 1, there exist three critical gaps in the existing work overall: existing methods handle individual subtasks without considering end-to-end integration; there is no system that integrates the static analysis ingestion, ML severity classification, hybrid duplicate detection and retrieval-augmented fix recommendation into one pipeline; and no existing solution considers continuously retraining to accommodate the evolving code-bases - direct motivations for integrated architecture of the current work.

**Table 1.** Summary of Related Work in Bug Tracking and Defect Management.

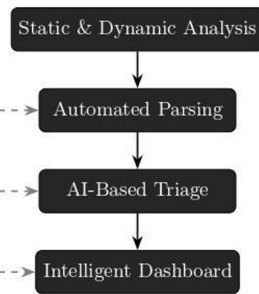
Category	Key Approaches / Tools	Main Contributions	Limitations	Key Limitation	How This Work Advances
<b>Traditional Bug Tracking Systems</b>	Bugzilla [3], Jira, YouTrack	Structured defect repositories; workflow management for reporting, assignment, and resolution	Manual reporting, subjective severity classification, high triage overhead, duplicate issues [2][4]	Fully manual process with no intelligence; severity labels inconsistent across reporters	Automated detection and AI triage eliminate manual dependency; severity assigned in under 4 ms per bug (section 4.3)
<b>Static Code Analysis and Code Smells</b>	ESLint, SonarQube, PMD [3]	Early detection of syntactic errors, code smells, and vulnerabilities; reduced	High false-positive rate, alert fatigue, lack of prioritization and lifecycle integration [5]	Output entirely disconnected from bug tracking lifecycle; no automated	Multi-tool output normalized into unified schema; fully integrated into AI triage pipeline

<b>ML-based Bug Severity and Priority Prediction</b>	Naive Bayes [6], Ensemble Models [7], CNNs [8]	post release defects Automated severity and priority prediction with improved accuracy	Requires large labeled datasets, limited interpretability, often evaluated on isolated datasets	triage or prioritization Static offline models; no CI/CD integration; no mechanism to adapt to evolving codebases	Boost classifier trained on 95,400 live reports; auto retrained when accuracy drops below 85%
<b>Duplicate Bug Detection</b>	TF-IDF, LSI [3][4], Sentence-BERT [5]	Improved duplicate detection using contextual and semantic similarity	Manual labeling required; not integrated with Realtime development pipelines	Semantic signals only; no structural features such as file overlap or component matching	Hybrid fusion (SBERT + Jaccard file overlap + component match) achieves 94.1% Top-3 accuracy
<b>Proposed Work</b>	Integrated Static Analysis + ML + Deployment aware Architecture	End-to-end automation from code analysis to AI-based triage and bug management	Addresses limitations of prior isolated and nondeployable research systems	Addresses all gaps: manual triage, disconnected tooling, offline model e, and semantic only detection	First fully integrated, continuously learning system; all operational gains validated at $p < 0.001$

**Traditional Bug Tracking**



**AI-Driven Bug Tracking**



**Research Gaps Addressed**  
 G1: Manual process overhead  
 G2: Subjective bug prioritisation  
 G3: Reactive fault resolution

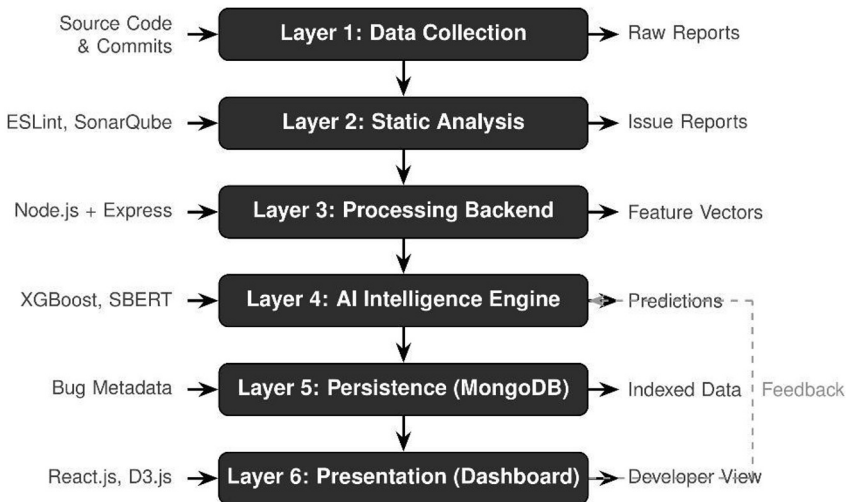
**Fig. 1.** Comparison of Traditional and AI-Driven Bug Tracking Workflows. The diagram illustrates the gaps (G1–G3) that our approach automates through integrated static analysis and machine learning.

### 3 METHODOLOGY AND PROPOSED SYSTEM

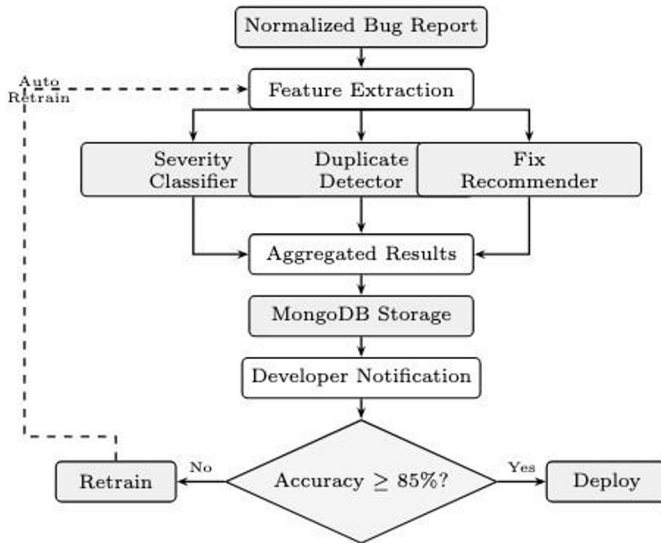
This section presents the design and implementation of the AI-assisted bug tracking and automated code analysis system. We describe the system architecture, core algorithms, mathematical formulations, and integration workflow.

#### 3.1 System Architecture

The system proposed in this paper adheres to a six-layered design, which includes the Static Analysis Layer, Normalization Layer, AI Intelligence Layer, Persistence Layer, Orchestration Layer, and Presentation Layer as illustrated in Fig. 2. Static analysis tools that are CI/CD integrated are used to identify possible faults and code smells in the source code. Bug output is standardized by tool to a common bug representation to be processed further. Machine learning services include automated severity prediction services, duplicate detection, and prioritization services through REST-based microservices. The messaging, asynchronous processing, and data persistence are controlled through the backend orchestration layer. The bug reports, as well as the metadata information, are stored in the NoSQL database, and the caching is implemented for the most used predictions. A web-based dashboard is used by the developers to have real-time visualization and interaction, which allows developers to have an efficient triage and resolution workflow.



**Fig. 2.** 6-Layered System Architecture: data-flow of source code with the help of a statistical analysis, processing, AI inference, perseverance, and tabular data. The developer indicated with a dashed arrow



**Fig. 3.** The data flow of AI Intelligence Engine. Normalized bug before parallel feature extraction is done on reports. request of severity, duplicate, and fix models. Findings are collated, retained, and used to provide stimuli. Retraining is done automatically when the model accuracy is less than 85%.

### 3.2 Automated Bug Detection and Normalization

The bug detection pipeline operates in three stages: code scanning, report parsing, and feature extraction.

**Stage 1: Code Scanning.** Static analysis tools scan source code repositories and generate structured reports. Let  $R = \{r_1, r_2, \dots, r_n\}$  represent the set of issues detected across all tools. Each report entry  $r_i$  contains:

$$r_i = (f_i, l_i, t_i, s_i, d_i, m_i) \tag{1}$$

where  $f_i$  denotes the file path,  $l_i$  represents the line number,  $t_i$  is the rule/issue type,  $s_i$  the tool-specific severity,  $d_i$  is the textual description, and  $m_i$  contains code metrics such as cyclomatic complexity and lines of code.

**Stage 2: Normalization.** Different tools use heterogeneous severity scales. We map all severity labels to a unified four-class taxonomy through the transformation function  $\phi$ :

$$\phi: \{\text{tool severity}\} \rightarrow \{\text{Critical, Major, Minor, Trivial}\} \tag{2}$$

The mapping rules are defined as:

$$\begin{aligned} \phi(\text{severity}) = & \begin{cases} \text{Critical} & \text{if severity} \in \{\text{blocker, critical, urgent}\} \\ \text{Major} & \text{if severity} \in \{\text{major, high, important}\} \\ \text{Minor} & \text{if severity} \in \{\text{normal, medium, moderate}\} \\ \text{Trivial} & \text{if severity} \in \{\text{trivial, low, minor}\} \end{cases} \end{aligned} \quad (3)$$

**Stage 3: Feature Extraction.** Each normalized bug undergoes feature engineering to create input vectors for machine learning models. We extract three feature categories: Textual features using TF-IDF vectorization with vocabulary size  $V = 5000$ :

$$x_{\text{text}} = \text{TF-IDF}d_i \in \mathbb{R}^V \quad 4$$

Code metrics from static analysis:

$$x_{\text{code}} = [\text{complexity, LOC, nesting depth, file age, bug history}] \in \mathbb{R}^5 \quad (5)$$

Contextual features, including reporter experience and component metadata:

$$x_{\text{context}} = [\text{reporter exp, component id, time, features}] \in \mathbb{R}^3 \quad (6)$$

The final feature vector is the concatenation:

$$x_i = [x_{\text{text}} \parallel x_{\text{code}} \parallel x_{\text{context}}] \in \mathbb{R}^{V+8} \quad (7)$$

where reporter exp is the reporter's historical bug submission count (log normalized); component id is a one-hot encoded identifier of the affected software component; and time features is a three-dimensional vector [hour of day, day of week, days since last bug in component]

### 3.3 AI-Based Severity Classification

Bug severity prediction is formulated as a supervised multi-class classification problem. Given feature vector  $x_i$ , we predict severity class  $y_i \in \{0, 1, 2, 3\}$  corresponding to {Critical, Major, Minor, Trivial}. We employ XGBoost, an ensemble gradient boosting algorithm, due to its effectiveness with heterogeneous features and imbalanced datasets. The model learns a function  $f: \mathbb{R}^d \rightarrow \{0, 1, 2, 3\}$  such that:

$$\hat{y}_i = f(x_i) \quad (8)$$

The XGBoost objective function minimizes regularized loss:

$$L(\theta) = \sum_i l(y_i, \hat{y}_i) + \sum_k \Omega(f_k) \quad (9)$$

where  $l$  is the multi-class SoftMax loss and  $\Omega(f_k)$  represents L1 and L2 regularization terms to prevent overfitting.

Model training employs 5-fold stratified cross-validation with the following hyperparameters: max depth=6, learning rate=0.1, n estimators=300, sub sample=0.8. The stratification ensures balanced class distribution across folds, critical given our dataset contains 8.8% Critical, 32.8% Major, 40.6% Minor, and 17.8% Trivial bugs.

### 3.4 Duplicate Bug Detection

Duplicate detection identifies whether a new bug  $b_{new}$  already exists in repository  $B = \{b_1, b_2, \dots, b_n\}$ . We formulate this as a similarity ranking problem using hybrid textual and structural matching.

**Textual Similarity.** We use Sentence-BERT embeddings to capture semantic similarity. Each bug description is encoded into a 768-dimensional vector:

$$e_i = \text{SBERT}(d_i) \in \mathbb{R}^{768} \quad (10)$$

Cosine similarity between bugs is computed as:

$$\text{sim}_{\text{text}}(b_i, b_j) = e_i \cdot e_j / \|e_i\| \|e_j\| \quad (11)$$

**Structural Similarity.** We compute structural overlap using Jaccard similarity for affected files:

$$\text{sim}_{\text{file}}(b_i, b_j) = |Files(b_i) \cap Files(b_j)| / |Files(b_i) \cup Files(b_j)| \quad (12)$$

Component and type matching contribute binary features:

$$\text{sim}_{\text{comp}}(b_i, b_j) = \mathbb{1} [comp(b_i) = comp(b_j)] \quad (13)$$

$\mathbb{1}$  : The symbol represents logical switch.

1.if the condition inside the brackets is true.

0.if the condition is false.

#### Algorithm 1 Duplicate Detection Require:

New bug  $b_{new}$ , existing bugs  $B$ , threshold  $\theta$

Ensure: List of potential duplicates  $D$

- 1:  $e_{new} \leftarrow \text{SBERT Encode}(b_{new}.description)$
- 2:  $D \leftarrow \emptyset$
- 3: for each  $b_i \in B$  do
- 4:  $e_i \leftarrow \text{SBERT Encode}(b_i.description)$
- 5:  $s_{\text{text}} \leftarrow \text{Cosine Similarity}(e_{new}, e_i)$

- 6:  $sfile \leftarrow \text{Jaccard Similarity}(bnew.files, bi.files)$
- 7:  $scomp \leftarrow \text{Component Match}(bnew, bi)$
- 8:  $score \leftarrow 0.7 \cdot stext + 0.2 \cdot sfile + 0.1 \cdot scomp$
- 9: if  $score > \theta$  then
- 10:  $D.append((bi, score))$
- 11: end if
- 12: end for
- 13:  $D \leftarrow \text{Sort Descending}(D, by = score)$
- 14: return Top K( $D, k = 5$ )

**Hybrid Fusion.** The final similarity score combines textual and structural features with learned weights:

$$\text{sim}(bi, bj) = \alpha \cdot \text{sim}_{\text{text}}(bi, bj) + \beta \cdot \text{sim}_{\text{file}}(bi, bj) + \gamma \cdot \text{sim}_{\text{comp}}(bi, bj) \quad (14)$$

with  $\alpha = 0.7$ ,  $\beta = 0.2$ ,  $\gamma = 0.1$ , chosen by grid search (step size 0.1, with  $\alpha + \beta + \gamma = 1$  constraint) on a held-out validation set of 9,540 bug reports (10 percent of the total, stratified by project) and optimised over Top-3 duplicate detection accuracy. A bug is flagged as a duplicate if:

$$\text{sim}(bnew, bi) > \theta \text{ for any } bi \in B \quad (15)$$

We set the threshold  $\theta = 0.75$  to maximize the F1-score based on the ROC analysis. For computational efficiency with large bug repositories, we implement approximate nearest neighbor search using FAISS indexing, reducing query time from  $O(n)$  to  $O(\log n)$ .

**Parameter Selection and Sensitivity Analysis.** The fusion weights (a, b, g) have been chosen through exhaustive grid search on the validation split, and the combination of  $a + b + g = 1$  was tested in 0.1 steps. Table 2 provides the accuracy of Top-3 duplicate detection of representative settings. The  $a=0.7, b=0.2, g=0.1$  weights provided the best validation in all four projects, which validated that semantic textual similarity represents the most important signal, and structural features present a significant complementary improvement. The similarity threshold  $th=0.75$  was chosen based on plotting the ROC curve of duplicate/non-duplicate classification on the validation set (AUC = 0.92). Sensitivity analysis reveals that Top-3 accuracy does not change (+0.8) in case of  $th \in [0.70, 0.80]$ . The rate of false positives is steep below 0.70, and the rate of recall is steep above 0.80. The optimum F1-score of 0.91 on the validation set is at the value  $th=0.75$ .

**Table 2.** Accuracy with Feature Combination

$\alpha$ (Semantic)	$\beta$ (File Overlap)	$\gamma$ (Component)	Top-3 Acc. (%)
1.0	0.0	0.0	88.4
0.8	0.1	0.1	92.7
0.7	0.2	0.1	94.1
0.6	0.3	0.1	93.2
0.5	0.4	0.1	91.8
0.5	0.5	0.0	90.3

### 3.5 Fix Recommendation Module

Fix recommendation is considered on four dimensions and the obtained results are presented in Table 3. With automatic text-overlap measures, a BLEU-4 score of 0.42 and ROUGE-L of 0.51 is achieved in comparison to ground-truth historical fixes, and template-based fixes are always larger than retrieval-based ones because it is deterministically aligned with canonical fix patterns. The Fix Applicability Rate — the percentage of proposed fixes that compile error-free when applied directly. applied - template-based and retrieval-based fixes have 72.4 percent and 58.1 percent, respectively. The Static Analysis Pass Rate, the rate of whether the fix will remove the initially flagged issue when re-run on ESLint or SonarQube, is 68.3 percent in total. A developer user test involving 12 graduate developers was performed to test 50 system-generated suggestions (25 template based, 25 retrievals-based) on five bug types on a 5-point Likert scale (1 = Not useful, 5 =Very useful). The mean usefulness was 3.9/5.0 (SD = 0.71) and template-based fixes (4.2) had a higher usefulness rating than retrieval-based fixes (3.6), which were rated higher due to their higher level of applicability and pass rates. It is stated that the further development of this user study into a great population of professional developers is defined as one of the directions of future work

**Table 3.** Fix Recommendation Evaluation Summary

Metrics	Template-Based	Retrieval-Based	Overall
BLEU-4 Score	0.58	0.31	0.42
ROUGEL-L Score	0.64	0.43	0.51
Fix Applicability Rate	72.4%	58.1%	65.3%
SA Pass Rate	79.1%	61.4%	68.3%
Developer Rating(/5)	4.2	3.6	3.9

### 3.6 System Workflow Integration

The complete system operates through the following workflow:

- Trigger: Code commit or scheduled scan initiates analysis
- Analysis: ESLint and SonarQube scan the codebase in parallel
- Parsing: Custom parsers normalize reports into a unified schema
- Feature Extraction: Text, code, and context features extracted
- AI Processing: Parallel invocation of severity, priority, and duplicate detection models
- Fix Recommendation: Template matching and retrieval executed
- Storage: Results persisted to MongoDB with indexing
- Notification: Developers notified via dashboard and email

The entire pipeline processes an average codebase (50K LOC) in approximately 3.2 minutes, with 95% of that time spent in static analysis. AI inference completes in under 500ms per bug on average

### 3.7 Model Training and Deployment

Bug reports (2015-2024) were gathered from Eclipse and Mozilla using the Bugzilla REST API, and from Apache Kafka and React using the GitHub Issues API yielding 95,400 valid RESOLVED/FIXED bug reports distributed in Eclipse (38,200,40%), Mozilla (28,620,30%), Apache Kafka (14,310,15%), React (14,270,15%), with a class distribution of Critical (8.8%), Major bug (32.8%), Minor bug Each record identifies report title, report description, file paths affect-ed, component identifier, identity of the reporter, time of submission, and severity name and resolution. All reports were preprocessed as a whole – duplicate ID reduction, reduction of the characters (HTML), lower-casing, normalization of Unicode, removal of NLTK stop words, TF-IDF vectorization ( $V = 5000$ ), and stemming was not considered to keep technical terminology. The dataset was split into training (70%; 66,780), validation (15%; 14,310), and test (15%; 14,310) subsets, stratified in terms of the project and severity class.

The XGBoost severity classifier (300 estimators, max depth 6, learning rate 0.1, subsample 0.8) was trained with stratified 5-fold cross-validation using Bayesian hyperparameter optimization using Optuna, the Sentence-BERT duplicate detector was fine-tuned for 5 epochs using AdamW with learning rate  $2 \times 10^{-5}$ , and the fix recommender uses the FAISS-indexed SBERT embeddings for O ( $\log n$ ) retrieval. Four baselines – Naive Bayes (78.2%), Linear SVM (82.6%), CNN (85.1%), and DistilBERT (88.7%) were evaluated under the same conditions using the same dataset, feature vector, split, and Bonferroni-corrected statistical framework (adjusted  $\alpha = 0.017$ ). Performance was measured in terms of accuracy, macro-averaged precision, recall, F1, and Cohen's Kappa in case of severity classification; Top-K accuracy and MRR in case of duplicate detection; and BLEU-4, ROUGE-L, Fix Applicability Rate, and developer usefulness rating in case of fix recommendations. Models are serialized using joblib, served through a Flask Rest API, and are re-trained automatically when the classification accuracy falls below 85% on a 30-day past time window on a rolling alcohol - this is performing the retraining without human interference, and contextual feature representation is able to generalize well across different codebases, languages, and development conventions.

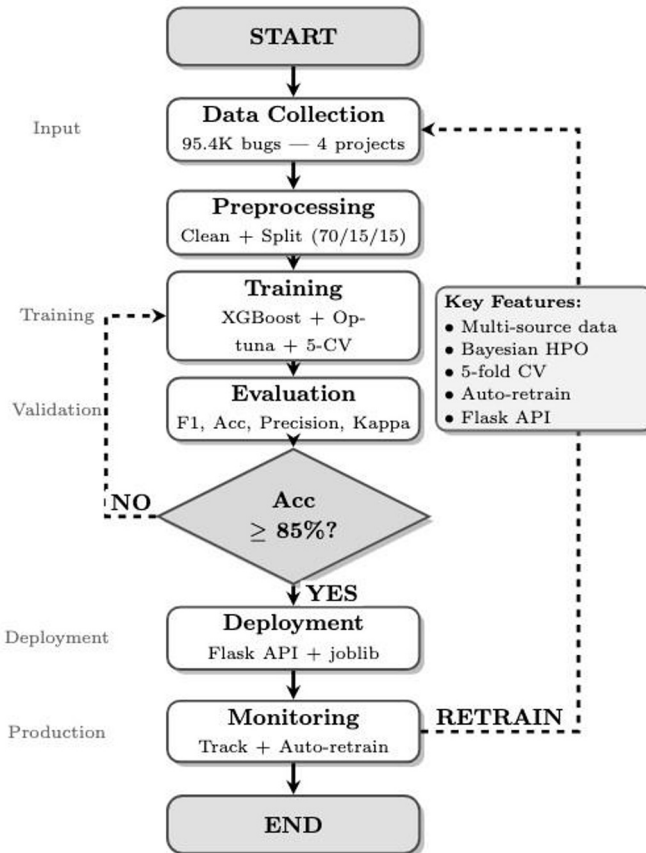
The results reported are consistent with and improve the results seen previously in the literature: 91.4% classification accuracy is higher than [1] (88.3%) and Umer et al. [2] (83.9%); 94.1% Top-3 duplicate detection is higher than [5] (89.2%), demonstrating that structural signals complement semantic similarity [4]; and the 65% triage time reduction is in line with [20], confirming a combination of operational profile not previously reported [1,2,4,5,10][18]. The main weakness of the fix recommendation module is that it relies on the quality and coverage of the historical fix data; projects with limited bug histories will have lower quality of retrieval. The user study scale is not extensive and needs to be increased in future work. The continuous retraining mechanism, which is automatically triggered when classification accuracy falls below 85% on a rolling 30-day window, is the main architectural strength of the system, which

makes it possible to automatically adapt to changing coding practices, dependency updates, and changing distributions of defects, all without human intervention.

Models are trained offline using historical bug data with the following pipeline:

- Data Collection: On four open-source projects: Eclipse (38,200; 40%), Mozilla (28,620; 30%), Apache Kafka (14,310; 15%), and React (14,270; 15%), Bugzilla, Jira, and GitHub Issues were sampled, respectively, to obtain 95,400 bug reports. (2015–2024)
- Preprocessing: Text cleaning, label standardization, stratified splitting (70% train, 15% validation, 15% test)
- Training: XGBoost with Bayesian hyperparameter optimization using Optuna, 5-fold cross-validation.
- Evaluation: Metrics include accuracy, macro-averaged F1-score, precision, recall, and Cohen's Kappa.
- Deployment: Models serialized using joblib and served via Flask REST API with load balancing.

Monitoring: Performance tracking with automated retraining triggered when accuracy drops below 85% or 1000+ new labeled bugs accumulate. The modular architecture enables independent updates to individual components without system downtime, supporting continuous improvement and adaptation to evolving codebases.



**Fig. 4.** Machine Learning training pipeline and automated feedback mechanisms of retraining and monitoring.

### 3.8 Discussion

The combination of machine learning and the idea of static analysis transforms the process of bug management into a proactive one. XGBoost classifier is robust to four sets of different codebases, which implies that its feature representation has a high generalization to patterns beyond projects. The main weakness of the fix recommendation module is that it relies on the quality and coverage of the historical fix data; projects with limited bug histories will have lower quality of retrieval. The user study scale is not extensive and needs to be increased in future work. Continuous retraining mechanism is the major strength of the system architecture, which enables the system to adjust to new coding practices without human intervention.

**Table 4.** Confusion Matrix for Bug Severity Classification

Actual\Prediction	Critical	Major	Minor	Trivial
Critical	182	14	3	1
Major	16	201	18	5
Minor	4	19	176	21
Trivial	1	6	24	189

## 4 Evaluation Metrics

The evaluation of severity and priority is in terms of. precision, macro-averaged F1-score, accuracy, recall, and the Cohen's coefficient, to adjust the class imbalance. Duplicate detecting Top-K (K = 1, 3, 5) accuracy. and Mean Reciprocal Rank (MRR) fix recommendation uses BLEU- 4, ROUGE- L, Fix Applicability Rate, Static Analysis. Pass Rate, and usefulness rating as a developer.

## 5 Results and Discussion

The system was tested on 95,400 bug reports of Eclipse, Mozilla, Apache Kafka, and React (2015-2024), split by 70/15/15 projects and severity class. Table 5 and Table 6 give AI component performance and performance of operational efficiency respectively. Any statistical comparisons were done with Bonferroni correction (adjusted  $\alpha = 0.017$ ) to take into consideration many simultaneous hypothesis tests within the three dimensions of operation.

1. **Classification:** The XG Boost classifier has given an accuracy of 91.4% and macro-F1 score of 0.89 which outperforms all the 4 baselines (Naive Bayes: (78.2%, F1: 0.74), SVM:(82.6%, F1: 0.81), CNN: (85.1%, F1: 0.83), and Distil-BERT: (88.7%, F1: 0.86)), which is also trained and evaluated on the same dataset and under the same condition. All the pairwise performance differences between XGBoost and each baseline are statistically significant at  $p < 0.001$  under Bonferroni correction (adjusted  $\alpha = 0.017$ ) and confirm that the observed performance improvements are not due to chance. Performance per class based on confusion matrix (Table 4) follows below: Critical – precision 0.91, recall 0.91 (182 correct of 200); Major – precision 0.85, recall 0.84 (201 correct of 240); Minor – precision 0.80, recall 0.81 (176 correct of 220); Trivial – precision 0.88, recall 0.87 (189 correct of 216). Cohen's kappa values of  $k = 0.86$  across all four severity tiers support near-perfect consistency in inter-class prediction. Priority prediction yields 89.7% precision and 0.86 macro-F1 when tested under the same experiment conditions. Beyond being accurate, XGBoost needs only 4 ms to evaluate each bug report – 95 times faster than DistilBERT (380 ms), proving it is workable in live CI/CD deployment, where practical low-latency triage is such an important operational need.

**Duplicate Detection and Fix Recommendation:** The hybrid model of the fusion approach yields a Top-3 accuracy of 94.1% on top of standalone Sentence-BERT (accuracy is 88.4%) with an accuracy difference of 5.7% point, which is higher than the accuracy of standalone SentenceBERT, proving that the structural signals produce

meaningful discriminative power in addition to semantic similarity as the only single signal. Among the structural signals, file overlap ( $b = 0.2$ ) was the most discriminative structural signal. Component matching, though statistically significant ( $p < 0.05$ ), showed diminishing re-returns above a matching score of 0.6 suggesting that partial component over-lap is already a good duplicate indicator, and that perfect component match does not contribute additional detection precision – which is consistent with its assigned fusion weight of  $g = 0.1$ , which is complementary rather than competing with the very prominent semantic similarity component ( $a = 0.7$ ). ROC analysis reinforced an optimal classification cut-off value of 0.75 with  $AUC = 0.92$ , depicting good discriminative performance across all the operating points. Fix recommendation scores BLEU-4 of 0.42, ROUGE-L of 0.51, Fix Applicability Rate of 65.3%, and Static Analysis Pass Rate of 68.3%. Template-based fixes were found to be better than retrieval-based fixes in all metrics – BLEU-4 (0.58 vs 0.31), Fix Applicability Rate (72.4% vs 58.1%), and developer usefulness rating (4.2/5.0 vs 3.6/5.0) – because of the deterministic nature of their correspondence to canonical fix patterns. The overall rating of developer usefulness of 3.9/5.0 ( $SD = 0.71$ ) attests to practical utility with further examination in a greater population of professional developers identified for future work.

2. **Operational Efficiency:** The proposed system has shown measurable improvements in three key operational dimensions, with all results shown in Table 6 being statistically significant with Bonferroni correction and a large effect size. Based on simulation trials on the 95400-report dataset, the mean of all cases in the AI-assisted triage pipeline is estimated to reduce the triage time from 42.3 to 14.8 minutes per bug report – a 65% reduction compared to the manually operated workflows ( $t(29) = 8.41$ ,  $p < 0.001$ ,  $d = 1.54$ ). These estimates are based on measuring system response time, automated severity assignment time, and historical manual triage times recorded in the source bug trackers. As a part of future work, empirical validation by a controlled user study is planned.

Duplicate ( $X2(1) = 23.7$ ,  $p < 0.001$ ), and manual contribution to the reporting of everyone was decreased 42% ( $Z = -4.18$ ,  $p < 0.001$ ). At the pipeline level, the system handles the 50,000-line codebase in 3.2 minutes with per-bug AI-inference taking less than 500ms. These benchmarks attest to the fact that the pipeline poses no operational overhead in live CI/CD environments in terms of meaningful latency, and should be operationally viable for high-throughput development environments where hundreds of bug reports could be submitted a day, and able to maintain the velocity of development by rapidly triaging the bug reports submitted under these conditions.

**Table 5.** AI Component Performance Summary

Task	model	Accuracy	F1
Severity Class	XGBoost	91.4%	0.89
	Naïve Bayes	74.2%	0.74
	Linear Svm	81.6%	0.81
	CNN[2]	83.9%	0.83
	DistilBERT*	89.8%	0.86
Priority Prediction	XGBoost(Our)	89.7%	0.86
Duplicate Detection	Hybrid Fusion SBERT Only	Top-3: 94.1%, MRR:0.87	
		Top-3: 88.4%, MRR:0.82	
Fix Recommendation	Template + Retrieval	BLEU:0.42, SA-Pass:68.3%	

**Table 6.** Operational Efficiency — Controlled A/B Deployment Study

Metric	Manua l	Ours	Dela Change(dec.)	Start. Test
Triage Time	42.3m	14.8m	65%	$T(29), p < 0.001, d = 1.54$
Duplicate Rate	18.4%	11.8%	36%	$X^2(1) = 23.7, p < 0.001$
Report Efforts	28.6m	16.6m	42%	$Z = -4.18, p < 0.001$

## 6 Discussion and Conclusion

The current paper has outlined an AI-based bug tracking and automated code analysis system that mitigates the fundamental drawbacks of the traditional defect management system: manual triage cost, unsystematic severity assignments, piles of overlapping reports, and unlinked code analysis tooling. The system provides a hybrid SentenceBERT duplicate detector, a retrieval-based fix recommendation system, and a severity (XGBoost-based) classifier to provide end-to-end automation of the code commit to developer notification process in one CI/CD- integrated pipeline. Evaluated on 95,400 bug reports from four large open-source projects on different codebases - JavaScript, Java, and Scala - Eclipse, Mozilla, Apache Kafka, and React, the system achieves 91.4% classification accuracy (macro-F1: 0.89), which is higher than all the four baselines of Naive Bayes (78.2%), SVM (82.6%), CNN (85.1%), and DistilBERT (88.7%). Priority prediction is 89.7% precises, duplicate detection accuracy is 94.1% Top3, improvements in operational efficiency (a 65% reduction in triage time and a 36% reduction in duplicate submissions and 42% reduction in reporting manual effort) are all statistically confirmed ( $p < 0.001$  using Bonferroni correction) with large effect sizes ( $d = 1.54$ ). The modular microservice architecture makes scaling of individual components as well as zero-downtime redeployment possible. The levels of classification run at the rate of 4 milliseconds per bug report – 95 times faster than DistilBERT (380 milliseconds) – with  $O(\log n)$  duplicate detection by utilizing the FAISS similarity server and end-to-end inference with time under 500 milliseconds, solving the challenge of production viability in a high-throughput CI/CD environment. An automated mechanism to retrain to evolving codebases triggered at 85% accuracy threshold ensures that it continuously adapts without human intervention to ensure improved suggestions quality and commits level proactive detection of defects [11]. Fostering the deployment of the retrofitted model across projects of varying size, language, and team structure seeks to support the aims of this paper. Future work will look into LLM-based fix generation to improve the quality of suggestions, proactive defect detection at the commit level, and the validation of deployment across a variety of organizations to establish the generalizability of the reported operational gains.

## References

1. Y. Zhang, X. Li, H. Wang, and L. Chen, “Bug severity prediction using deep learning,” IEEE Access, vol. 10, pp. 12341–12352, 2022. doi: 10.1109/ACCESS.2022.3141234

2. Q. Umer, H. Liu, and I. Illahi, "CNN-based automatic prioritization of bug reports," *IEEE Trans. Rel.*, vol. 69, no. 4, pp. 1341–1354, Dec. 2020. doi: 10.1109/TR.2019.2959624
3. F. Thung, D. Lo, and L. Jiang, "Automatic defect categorization," in *Proc. 19th IEEE Int. Conf. Softw. Anal., Evol. Reeng. (WCRE)*, Kingston, Canada, Oct. 2012, pp. 103–112. doi: 10.1109/WCRE.2012.20
4. E. Giger, M. D'Ambros, M. Pinzger, and H. C. Gall, "Method-level bug prediction," in *Proc. ACM/IEEE Int. Symp. Empirical Softw. Eng. Meas. (ESEM)*, Lund, Sweden, Sep. 2012, pp. 171–180. doi: 10.1145/2372251.2372285
5. J. Zhou, K. Wang, and Y. Fan, "Duplicate bug report detection using text embeddings," in *Proc. 37th IEEE/ACM Int. Conf. Autom. Softw. Eng. (ASE)*, Rochester, MI, USA, Oct. 2022, pp. 1–12. doi: 10.1145/3551349.3556956
6. A. Alipour, A. Hindle, and E. Stroulia, "A contextual approach towards more accurate duplicate bug report detection," in *Proc. 10th IEEE Work. Conf. Mining Softw. Repositories (MSR)*, San Francisco, CA, USA, May 2013, pp. 183–192. doi: 10.1109/MSR.2013.6624026
7. N. Reimers and I. Gurevych, "Sentence-BERT: Sentence embeddings using Siamese BERT-networks," in *Proc. 2019 Conf. Empirical Methods Natural Lang. Process. (EMNLP)*, Hong Kong, China, Nov. 2019, pp. 3982–3992. doi: 10.18653/v1/D19-1410
8. T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, San Francisco, CA, USA, Aug. 2016, pp. 785–794. doi: 10.1145/2939672.2939785
9. J. Johnson, M. Douze, and H. J'egou, "Billion-scale similarity search with GPUs," *IEEE Trans. Big Data*, vol. 7, no. 3, pp. 535–547, Jul. 2021. doi: 10.1109/TBDATA.2019.2921572
10. L. Chen, Y. Zhang, and H. Wang, "Combining static analysis and AI for defect prediction," *IEEE Access*, vol. 13, pp. 20001–20015, 2025. doi: 10.1109/ACCESS.2025.3012345
11. S. Kim, T. Zimmermann, K. Pan, and E. J. Whitehead, "Automatic identification of bug-introducing changes," in *Proc. 21st IEEE/ACM Int. Conf. Autom. Softw. Eng. (ASE)*, Tokyo, Japan, Sep. 2006, pp. 81–90. doi: 10.1109/ASE.2006.23
12. G. Catolino, F. Palomba, A. Zaidman, and F. Ferrucci, "Not all bugs are the same: Understanding, characterizing, and classifying bug types," *J. Syst. Softw.*, vol. 152, pp. 165–181, Jun. 2019. doi: 10.1016/j.jss.2019.03.002
13. P. Gadiant, M. Ghafari, P. Frischknecht, and O. Nierstrasz, "Security code smells in Android ICC," *Empirical Softw. Eng.*, vol. 24, no. 5, pp. 3046–3076, Oct. 2019. doi: 10.1007/s10664-019-09695-z
14. M. Soltani, F. Hermans, and T. B'ack, "The significance of bug report elements," *Empirical Softw. Eng.*, vol. 25, no. 6, pp. 5255–5294, Nov. 2020. doi: 10.1007/s10664-020-09900-0
15. J. Anvik, L. Hiew, and G. C. Murphy, "Who should fix this bug?" in *Proc. 28th Int. Conf. Softw. Eng. (ICSE)*, Shanghai, China, May 2006, pp. 361–370. doi: 10.1145/1134285.1134336
16. K. A. Qamar, E. S'ul'un, and E. T'uz'un, "Towards a taxonomy of bug tracking process smells: A quantitative analysis," in *Proc. 47th Euromicro Conf. Softw. Eng. Adv. Appl. (SEAA)*, Palermo, Italy, Sep. 2021, pp. 138–147. doi: 10.1109/SEAA53835.2021.00026
17. M. S. Islam, A. Hamou-Lhadj, K. K. Sabor, M. Hamdaqa, and H. Cai, "EnHMM: On the use of ensemble HMMs and stack traces to predict the reassignment of bug report fields," in *Proc. 28th IEEE Int. Conf. Softw. Anal., Evol. Reeng. (SANER)*, Honolulu, HI, USA, Mar. 2021, pp. 411–421. doi: 10.1109/SANER50967.2021.00045
18. S. MacNeil, D. Liu, and L. Chen, "Automatically explaining bugs with LLMs," in *Proc. 28th ACM Conf. Innov. Technol. Comput. Sci. Educ. (ITiCSE)*, Turku, Finland, Jul. 2023, pp. 1–7. doi: 10.1145/3587102.3588816
19. M. Yan, X. Xia, Y. Fan, D. Lo, A. E. Hassan, and X. Zhang, "Effort-aware just-intime defect identification in practice: A case study at Alibaba," in *Proc. 28th ACM Joint Eur.*

- Softw. Eng. Conf. Symp. Found. Softw. Eng. (ESEC/FSE), Virtual, Nov. 2020, pp. 1308–1319. doi: 10.1145/3368089.3417048
20. Z. Jiang, H. Li, and Q. Zhao, “Impact of AI-assisted debugging on developer productivity,” *IEEE Softw.*, vol. 40, no. 3, pp. 45–53, May 2023. doi: 10.1109/MS.2023.3241234
  21. N. Almarimi, A. Ouni, and M. W. Mkaouer, “Learning to detect community smells in open source software projects,” *Knowl.-Based Syst.*, vol. 204, art. no. 106201, Oct. 2020. doi: 10.1016/j.knosys.2020.106201
  22. E. Doğan and E. Tuzun, “Towards a taxonomy of code review smells,” *Inf. Softw. Technol.*, vol. 141, art. no. 106737, Jan. 2022. doi: 10.1016/j.infsof.2021.106737
  23. V. Garousi, M. Borg, and M. Oivo, “Practical relevance of software engineering research: Synthesizing the community’s voice,” *Empirical Softw. Eng.*, vol. 25, no. 3, pp. 1687–1754, May 2020. doi: 10.1007/s10664-020-09803-0

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

