



SCUOLA
ALTI STUDI
LUCCA

Scuola IMT Alti Studi Lucca

Will It Break in Production? Metric-Driven Prediction of Residual Defects in Python Systems

Questa è la versione preprint della seguente opera:

Original

Will It Break in Production? Metric-Driven Prediction of Residual Defects in Python Systems / De Rosa, G., Liguori, P.. - (2026).

Availability:

This version is available at: 20.500.11771/41660

Publisher:

Published

DOI:

Terms of use:

This publication is made accessible in accordance with the terms for deposit in the institutional repository, as defined by the IMT School for Advanced Studies Lucca's Open Access Policy. (https://library.imtlucca.it/sites/default/files/regolamento-policy-open-access-imtlib_0.pdf).

Si prega di consultare le pagine informative dell'editore relative alle politiche di autoarchiviazione.

(Article begins on next page)

Will It Break in Production? Metric-Driven Prediction of Residual Defects in Python Systems

Giuseppe De Rosa^{*†}, Pietro Liguori^{*}

^{*}University of Naples Federico II, Naples, Italy

giuseppe.derosa20@unina.it, pietro.liguori@unina.it

[†]IMT School for Advanced Studies Lucca, Lucca, Italy

giuseppe.derosa@imtlucca.it

Abstract—Python’s dynamic nature complicates testing and increases the possibility that some defects evade detection, so an effective fault prediction becomes essential. We examine whether post-release faults can be predicted using modern ML and DL. Using a balanced dataset of over 4,000 labeled faults with 83 product, process, statistical, and Python-specific metrics plus normalized code representations, we conduct cross-project experiments. LLMs and unsupervised models fail to distinguish residual from non-residual faults, while supervised metric-based models (RandomForest, XGBoost, CatBoost) perform far better, yielding a 0.85-0.9 recall and cutting false negatives by an order of magnitude. Process metrics, especially age, churn, and developer-activity, alongside class and file size, consistently prove most predictive. Notably, the Principal Component Analysis shows that metrics and code embeddings occupy distinct regions of the representation space, suggesting that they capture complementary rather than redundant information.

Index Terms—Software Defect Prediction, Residual, Post-Release, Pre-Release, Python, Machine learning, Software Metrics

I. INTRODUCTION

Software failures continue to impose severe operational and economic consequences across modern digital infrastructures. From the \$460 million trading loss at Knight Capital [1] to large-scale transportation disruptions and safety-critical incidents [2], [3], the inability to prevent defects from reaching deployment environments remains a central challenge for software dependability. Industry analyses echo this urgency: according to Gartner, organizations that invest in “digital immunity”, i.e., the systematic detection, containment, and mitigation of defects, can reduce downtime by up to 80%, significantly improving resilience and service continuity [4]. Accurately identifying fault-prone components before release, therefore, plays a decisive role in reducing maintenance effort, improving testing strategies, and increasing trust in the ever-growing ecosystem of software-driven systems.

Despite substantial advances in defect prediction, one category of faults continues to elude even rigorous testing pipelines: *residual faults*, i.e., defects that remain unnoticed until the software is already in production. In the reliability literature, these are also referred to as *post-release faults*, since the defining characteristic of a residual defect is precisely that it survives all pre-release verification activities and manifests only after deployment [5], [6]. Using post-release failures as an operational definition is, therefore, standard practice: a defect

revealed in the field is, by construction, one that the testing and quality-assurance processes failed to expose.

These faults are uniquely problematic. They often surface only under realistic operational conditions that are difficult or impossible to reproduce in controlled test suites; they tend to arise from subtle interactions, rare execution paths, or latent corner cases; and, critically, they rarely exhibit clear syntactic or structural patterns that distinguish them from non-residual defects. This makes predicting which faults will escape pre-release detection intrinsically more challenging than predicting fault-proneness in general.

Recent research in machine learning (ML) and deep learning (DL) for defect prediction has shown promising results by leveraging combinations of product and process metrics, along with learned code representations [7]. Yet several longstanding issues remain unresolved, and they become even more critical when the target is residual faults rather than generic defect-proneness. First, cross-project generalization is notoriously weak: models trained on one project frequently underperform when applied to others [8]. This limitation is exacerbated by the rarity of residual faults and their project-specific characteristics. Second, DL approaches often lack interpretability, limiting their usefulness in safety-critical settings where understanding why a component is predicted as risky is as important as the prediction itself [7], [9]. Third, residual faults are inherently rare, creating a challenging class imbalance that distorts learning and inflates apparent performance [10], [11]. These limitations motivate the need for carefully designed datasets and a broad spectrum of software metrics to capture language-specific behaviors that may influence residual fault manifestation.

This paper investigates residual fault prediction. Predicting residual faults is not equivalent to standard defect prediction: traditional defect predictors estimate whether a defect exists in a component, implicitly conflating two statistically different populations: faults that will be caught during testing and faults that will escape into production. We formulate a distinct prediction target: estimating the probability that a defect, once introduced, survives all pre-release validation activities and reaches production. This distinction is operationally significant because it enables targeted interventions such as selective manual review or fuzzing only for components with high escape risk rather than the entire codebase. To support this

task, we leverage a diverse range of traditional software metrics, including product metrics that quantify size, control-flow complexity, documentation quality, coupling, and inheritance structures; Python-specific syntactic indicators; statistical dispersion measures; and process metrics that capture code churn, temporal evolution, commit activity, and developer behavior. We focus on Python because of its widespread adoption in modern software infrastructures and its dynamic, interpreted nature, which makes defects more likely to manifest only at runtime [12]. These characteristics increase the probability that subtle faults escape pre-release testing, making Python a particularly relevant and challenging setting for studying residual defects [13], [14].

Our study yields three key contributions:

- 1) **A new dataset of residual and non-residual Python faults.** We build a balanced and reproducible dataset of more than 4,000 classified faults from real-world Python projects, complemented by 500 additional cases from BugsInPy. Each instance includes code, software metrics, and version-control information, enabling systematic experimentation on residual-fault prediction and supporting a repeatable evaluation.
- 2) **A systematic assessment of LLM-based code representations for residual fault prediction.** We evaluate both proprietary and fine-tuned open LLMs under uniform conditions. Although their embeddings contain meaningful signals, the models themselves perform unreliably: closed-source systems achieve F1-scores between 0.43 and 0.72, often collapsing into degenerate behavior, while fine-tuned models reach only 0.35–0.38 and yield a recall of 0.3 on average. A representation-space analysis shows that embeddings remain largely orthogonal to metrics (maximum correlation 0.27; mean 0.24), indicating that LLMs encode complementary but insufficient information for this task.
- 3) **An investigation of software metrics through lightweight ML models.** We assess how far one can go using only traditional software metrics and simple numerical models. Unsupervised anomaly detectors perform poorly, confirming that residual faults are not statistical outliers. In contrast, supervised learners such as RandomForest, XGBoost, and CatBoost achieve F1-scores of 0.71–0.73 and reach 0.85–0.9 recall, reducing false negatives by an order of magnitude compared to LLMs. Feature-importance and SHAP (i.e., a method that assigns each feature a quantitative contribution to the model’s prediction) analyses highlight process, size, and developer-activity metrics as the most informative families, revealing stable structural and historical patterns that support effective residual-fault prediction.

To foster both academic and industrial research on this topic, we publicly release our replication package [15], which includes the datasets and scripts used for the analyses in the paper.

The remainder of this paper is organized as follows. Sec-

tion II introduces the background and motivates our focus on residual faults. Section III details the data collection, metric extraction, and model design. Section IV describes the experimental setup, while Section V presents our evaluation results. Section VI discusses the results. Related work is discussed in Section VII, and threats to validity are outlined in Section VIII. Finally, Section IX concludes with reflections and directions for future work.

II. BACKGROUND

Software defect prediction aims to proactively identify modules within a system that are likely to contain undiscovered bugs before these defects manifest during real-world operation. Typically, this involves collecting diverse code metrics, such as complexity, size, object-oriented design indicators, change history, and developer-related factors, and applying machine learning algorithms to detect patterns associated with defect-prone components [16]. Conventional approaches train classification models on historical defect data, allowing development teams to prioritize testing and debugging efforts [17].

However, most existing prediction techniques focus primarily on defects identified during pre-release testing, neglecting those that surface only after deployment. These post-release issues, commonly referred to as *residual defects*, exhibit fundamentally different characteristics. They often affect components that appeared stable during testing, but fail in production under atypical inputs, varied deployment contexts, or unanticipated runtime interactions. Because of their contextual and environment-dependent nature, residual defects tend to be harder to detect, more costly to resolve, and more disruptive when uncovered in operational settings [6].

Residual defects are especially concerning in systems developed using dynamically typed languages. Among them, Python has experienced explosive growth across a wide range of domains, becoming foundational in web development (e.g., Django, Flask), machine learning (e.g., TensorFlow, scikit-learn), and data analytics (e.g., Pandas, NumPy [18], [19]). While its flexibility and expressive syntax contribute to its popularity, Python’s dynamic typing and interpreted execution model make it particularly susceptible to runtime errors, such as type mismatches, undefined variable access, and improper exception handling, that are difficult to detect via static analysis or conventional testing [20]. This susceptibility introduces significant operational risk, especially in large-scale and distributed environments [21].

III. METHODOLOGY

This section describes the methodology used to predict residual faults. Fig. 1 presents an overview of the proposed workflow, which consists of four main stages: (i) data collection from PyResBugs and real-world Python projects; (ii) classification of each fault as residual or non-residual to build the training and test sets; (iii) collection of product, process, and statistical software metrics for every instance; and (iv) training and evaluation of both ML and DL models, including preprocessing and explainability analyses.

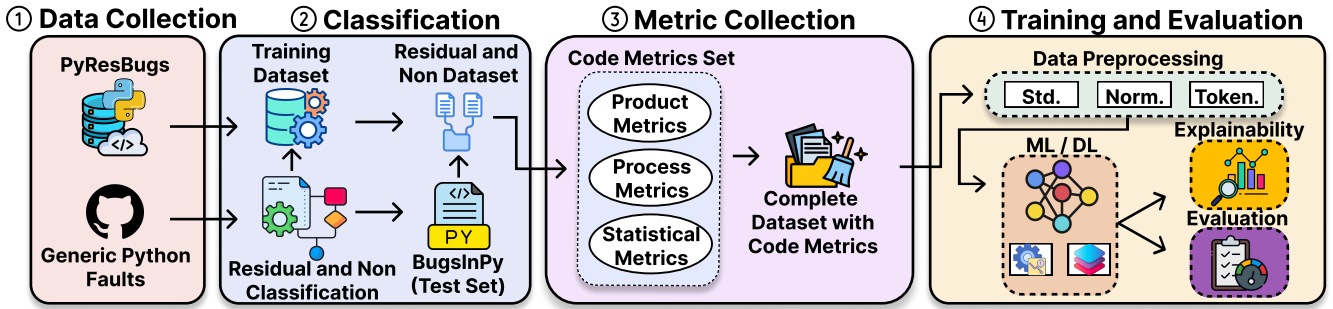


Fig. 1: Detailed methodology adopted in this work.

A. Data Collection

The success of fault prediction models largely depends on the quality, representativeness, and balance of the underlying datasets. Our data collection process carefully addresses these concerns, systematically gathering examples of both residual (post-release) and non-residual (pre-release) faults.

For residual faults, we leveraged the publicly available *PyresBugs* dataset, which contains 5,007 documented post-release bugs extracted from 76 prominent Python frameworks and libraries [22]. Residual faults are bugs that escape pre-release testing and only surface after the software is deployed. *PyresBugs* was initially constructed by aggregating faults from two primary sources. First, faults were curated from eleven existing datasets published in major software engineering conferences and journals over the past five years. This step ensured broad coverage of well-established, community-validated fault examples. Second, the dataset was augmented with additional categories of residual faults, specifically targeting those challenging to identify in standard testing phases, such as time-related issues and concurrency errors. To reliably identify genuine residual faults, we employed a multi-stage filtering procedure: bug reports from software repositories (e.g., GitHub) were analyzed to identify defects explicitly discovered post-release. Each identified defect was confirmed by examining associated commits to verify the availability of code-level fixes, enabling the extraction of both faulty and corrected code versions. Furthermore, we verified that each fault originated from source-code defects (within method-level), thus excluding cases involving documentation issues, configuration errors, or build system problems.

To construct the set of non-residual faults, we mined a large number of Python projects selected to resemble those included in *PyresBugs* in terms of size, maturity, and development activity. This alignment is essential to avoid systematic differences between the two defect categories that might arise from project-specific conventions. For each repository, we traversed the entire commit history and applied a keyword-based heuristic to identify commits related to defect correction. Following established practice in mining software repositories, we flagged a commit as bug-related if its message contained terms such as fix, fault, bug, defect, crash, issue, and so on. Keyword-based identification of bug-fixing commits is a standard approach in large-scale repository mining, and

although imperfect, it offers a good balance between coverage and precision when the goal is to gather a broad sample of natural faults [23], [24]. All bug-fixing commits discovered in this way were retained, and each associated defect instance was marked as a possible buggy commit.

The purpose of this collection was to gather a broad and diverse set of naturally occurring bugs. For each qualifying commit, we recorded project metadata, commit ID, file paths, the corresponding diff patches, the corresponding faulty version, and the issue or pull request ID. These patches provide sufficient information to characterize the fault and compute all product, process, and statistical metrics used in this study. The resulting collection forms a large corpus of authentic, possibly buggy commits. The next step will be to clean and classify them as non-residual faults by inspecting the associated issue.

B. Residual or Non-Residual Classification Task

To effectively perform this classification step, we rely on a dedicated heuristic that infers whether each bug-fixing commit resolves a pre-release or post-release fault, using only locally available historical repository information. The procedure `CLASSIFYCOMMIT` detailed in Algorithm 1 operates in stages, progressively incorporating increasingly weaker forms of evidence. The algorithm first checks whether the commit is linked to an issue (line 2). If so, it extracts decisive temporal signals from the issue creation time (line 3). Such signals include (i) whether the issue was created before the project’s first stable release, which unambiguously identifies it as a pre-release defect, and (ii) whether the issue contains explicit references to behaviour observed in a shipped version, like numeric version identifiers (e.g., v1.X, release X, build X), "affects-X.Y" labels, or regression tags, which indicate that the problem occurred in code already distributed to users. Issues created before the first stable release are immediately classified as pre-release (line 4), whereas those containing evidence tied to a released version are immediately classified as post-release (line 6), without considering weaker heuristics.

When no such hard evidence is available, the algorithm computes two soft scores, *pre_score* and *post_score*, each aggregating multiple weak heuristics (line 8). These include milestone information (closed milestones contributing to the post-release score, pre-release qualifiers contributing

Algorithm 1 Classify a commit as residual or non-residual

```
1: procedure CLASSIFYCOMMIT( $r$ )
2:   if linked issue exists then
3:     collect strong pre/post hints from issue time
4:     if some strong pre hint then return pre-release
5:     end if
6:     if some strong post hint then return post-release
7:     end if
8:     compute soft pre/post scores
9:     if  $pre\_score > post\_score$  then return pre-release
10:    else if  $post\_score > pre\_score$  then return post-
    release
11:    end if
12:    if author is external non-contributor then return
    post-release
13:    else if author is internal then return pre-release
14:    end if
15:    end if
16:    ( $cd, frd$ )  $\leftarrow$  commit date and first stable release date
    of repo
17:    if  $cd \neq \perp$  and ( $frd = \perp$  or  $cd < frd$ ) then return
    pre-release
18:    end if
19:    return unknown
20: end procedure
```

to the pre-release score), textual indicators in the issue body (internal-test markers increasing the pre-release score, bug-report templates and reproduction steps increasing the post-release score), and the reporter’s role (external reporters incrementing the post-release score, contributors incrementing the pre-release score). All soft indicators are boolean and combined through an unweighted sum, with no tuned parameters or learned weights, to ensure reproducibility across projects and to avoid project-specific bias. The full list of indicators is provided in our replication package [15].

After accumulating the scores, the algorithm compares the two totals. Suppose pre_score exceeds $post_score$. In that case, the majority of temporal and contextual evidence associated with the issue points to activities that typically occur *before* a release, e.g., developers discovering a defect during testing, code review, or stabilization. In other words, a higher pre_score suggests that the bug was fixed before deployment, and the commit is therefore classified as a pre-release fix. Conversely, a higher $post_score$ indicates that the evidence aligns more with characteristics of field failures (e.g., reports submitted after deployment, user-facing symptoms, production traces), supporting a post-release classification.

Each score is, by design, an imperfect proxy: for instance, an external reporter may occasionally file an issue discovered through source inspection rather than through use of a deployed version, and textual markers such as reproduction steps can appear in pre-release bug reports as well. The algorithm mitigates this by never relying on a single indi-

cator in isolation; instead, it aggregates all available signals and decides only when the balance of evidence favours one category. When the evidence is insufficient and/or conflicting, the commit is conservatively labeled *unknown* and excluded from subsequent analyses. When both scores remain tied, the algorithm falls back on the identity of the issue author (line 12) as a tie-breaker. We noticed that issues reported by external users are far more likely to reflect failures observed in the field: they interact only with deployed versions of the software; therefore, issues reported by non-contributors almost always correspond to failures observed in production. Instead, issues filed by contributors typically originate from internal testing or development before a release. The reporter identity, thus, represents a weak but meaningful signal for distinguishing between defects discovered during pre-release activities and those surfaced only after deployment. This signal is particularly valuable in cases where the temporal evidence extracted from timestamps is insufficient to disambiguate the classification (i.e., when pre_score and $post_score$ are tied). We use the identity of the *issue author* to determine internal vs. external origin: reporters with no prior contributions are treated as external, while reporters with commits or reviews are considered internal.

When no linked issue exists, the algorithm falls back to a timestamp-only rule (line 17): commits dated before the first stable release are labeled pre-release, while all other commits receive the *unknown* label. This could appear as an asymmetry: pre-release can be inferred from timestamps alone, but post-release cannot, because a commit made after a first release may still address a pre-release fault whose fix was delayed. However, it is well-known that faults escaping the testing phase and manifesting in production are far less than the ones surfacing in the testing phase, especially in mature projects with good test suites. As a consequence, the heuristic adopts a more conservative approach with the post-release faults to avoid introducing false positives, which we consider the more damaging threat to validity to the dataset quality.

To assess the accuracy of the heuristic, two authors independently inspected a random sample of 400 classified commits. Each commit was judged against its linked issue, commit message, and surrounding development context. The two raters achieved a Cohen’s $\kappa > 0.80$, indicating strong inter-rater agreement, and the heuristic’s labels matched the consensus ground truth in approximately 95% of cases.

When no linked issue exists, the algorithm resorts to commit-level timing (line 16), which is the only remaining temporal information available in the absence of a linked issue. It retrieves the commit date and the repository’s first stable release date. If the commit precedes the first stable release or the project has no stable release at all, the commit is classified as pre-release (line 17). This follows the standard assumption that post-release failures occur after the first official stable release. Moreover, the conservative nature of our heuristic ensures minimal risk of false post-release classifications.

We intentionally adopt a conservative policy and never classify a commit as post-release based solely on its timestamp.

TABLE I: Software Metrics Adopted

Category	Metric [Acronym]
Product Metrics	
Complexity	Method: Cyclomatic Complexity [CC], Maximum Nesting Depth [MND], Number of Paths [NP], <i>Halstead:</i> Difficulty [HD], Length [HL], Vocabulary [HV], Volume [HVOL], Effort [HEFF], Maintainability Index [HMI], Distinct Operators [HDOP], Distinct Operands [HDND], Total Operators [HTOP], Total Operands [HTOA] File: Cyclomatic Complexity [F-CC], Maximum Nesting Depth [F-MND], Number of Paths (log scale) [F-NPLOG]
Size	Method: Lines of Code [LOC], Blank Lines [BLOC], Declaration Lines [DLOC], Executable Lines [ELOC], Statement Count [STMT], Declaration Statements [DSTMT], Executable Statements [ESTMT], Number of Inputs [NIN], Number of Outputs [NOUT], Number of Exits [NE], Number of Early Exits [NEE] Class: Class Lines [CLLOC], Code Lines [CCODE], Declaration Lines [CDLOC], Executable Lines [CELOC], Number of Methods [NOM], Number of Declared Methods [NOM-A] Number of Instance Variables [NIV] File: Total Lines [F-TLOC], Code Lines [F-CLOC], Blank Lines [F-BLOC], Statement Count [F-STMT], Declaration Statements [F-DSTMT], Executable Statements [F-ESTMT]
Documentation	Method: Comment Lines [COMLOC], Comment to Code Ratio [CCR], Comment Lines w/ Preceding Code [CLWB], Comment Ratio (Preceding) [CCR-B] Class: Comment Lines [CCOM], Comment to Code Ratio [CCR-C] File: Comment Lines [F-COMLOC], Comment-to-Code Ratio [F-CCR]
Coupling & Inheritance	Method: Fan-In [FI], Fan-Out [FO], Component Reuse [CR] Class: Depth of Inheritance Tree [DIT], Base Classes [BCs], Derived Classes [DCs]
Python Specific	Maximum Indentation [PMI], Magic Numbers [PMN]
Statistical Metrics	
Distribution	Entropy [ENT]
Process Metrics	
Temporal	Age in Days [AGE], Bug Density [BD], Fix Commits [FC]
Churn-Based	Average Code Churn [ACCH], Maximum Code Churn [MCCH], Total Code Churn [TCCH], Total Modified Statements [TMS]
Commit-Based	Total Commits [TC], Commits Count [CMC], Maximum Commit LOC Changes [MCLC], Average Commit LOC Changes [ACLC], Total Commit Churn [TCC], Commit Churn Added [CCA], Commit Churn Deleted [CCD], Change Pattern Count [CPC]
Method-Based	Method Churn Added [MCA], Method Churn Deleted [MCD], Total Method Churn [TMC], Average Method LOC Changes [AMLC], Maximum Method LOC Changes [MMLC]
Developer-Based	Distinct Authors [DA], Average Developer Experience [ADE], Distinct Commit Names [DCN], Average Commits per Author [ACA], Average Code Churn per Author [ACCA]

Due to non-linear release workflows, maintenance branches, and backports, timestamp-only evidence is unreliable for inferring post-release origin. In these cases, we prefer the neutral label *unknown*. If none of the above conditions apply, meaning the algorithm cannot derive reliable temporal or role-based evidence, the commit is left as *unknown* (line 19).

To assess the reliability of this heuristic, we conducted a manual inspection of 400 randomly sampled software faults drawn from the collected commits. The sample was drawn uniformly across all projects, and each commit was manually labeled based on its complete issue and commit context. Manual annotation was performed by one author and cross-checked by a second author in ambiguous cases, until consensus was reached. Both authors involved in the process (a Ph.D. student and an Assistant Professor in Computer Engineering) have a solid background in software testing and fault classification. Our manual analysis revealed that $\sim 95\%$ of these cases were correctly classified as pre-release or post-release. The remaining instances were almost exclusively commits that were either

not true bugs (e.g., documentation updates, refactorings) or lacked any associated issue or temporal information, and were therefore legitimately labeled as *unknown* by the algorithm.

We remark that our goal with this heuristic is not to build a perfect classifier, but to obtain sufficiently reliable labels to support downstream modeling. The conservative design, combined with manual validation, suggests that residual mislabeling noise is limited and does not compromise the trends observed in our empirical analysis.

C. Metric Collection

A critical step in developing effective fault prediction models is selecting a comprehensive subset of software metrics that accurately characterizes residual and non-residual faults.

We began by conducting a thorough review of existing literature to collect metrics previously demonstrated as effective predictors of software faults [25]. This resulted in a broad pool of 83 metrics, categorized into three main groups: product metrics, which describe static characteristics of the source code, statistical metrics, which capture the differences in the

code variability, and process metrics, which capture historical and evolutionary aspects of the codebase (see Table I).

The metrics adopted in this study cover a wide range of code characteristics at the method, class, and file levels. We collected them by using Scitools Understand [26] and ad-hoc scripts. *Product metrics* describe static properties of the code and are grouped into complexity, size, documentation, coupling, inheritance, and Python-specific dimensions. Complexity indicators include Cyclomatic Complexity [CC], Maximum Nesting Depth [MND], and the Number of Paths [NP] at method level, with corresponding file-level variants such as [F-CC], [F-MND], and [F-NPLOG]. In addition, an extensive set of Halstead metrics (Difficulty, Length, Vocabulary, Volume, Effort, Maintainability Index, and operator/operand counts) captures the informational structure of the code and the cognitive effort required to understand it. Size metrics quantify dimensional properties at multiple granularities, including Lines of Code, declaration and executable statement counts, inputs and outputs, and the number of exits or early exits in methods. Class-level size measures (e.g., [CLLOC], [CDLOC], [NOM]) reflect structural regularities that influence maintainability. File-level size indicators complement these with global views of the code layout. Documentation-related metrics, including comment lines, comment-to-code ratios, and comments appearing before code, characterize the density and distribution of documentation across methods, classes, and files. Coupling and inheritance measures, such as Fan-In and Fan-Out at the method level, and the Depth of Inheritance Tree and the number of base or derived classes at the class level, provide information on dependency structure and the potential for fault propagation. Python-specific indicators, including maximum indentation and magic numbers, capture aspects that often correlate with readability and error-proneness.

Statistical metrics complement them with a distribution-based measure of code regularity. Entropy [ENT] was computed using KenLM, trained on the ETH Py150 corpus of 150k Python files [27]. After training the language model, the entropy of each function in our dataset was computed, allowing us to estimate how typical or atypical a piece of code is relative to a large and diverse Python codebase.

Process metrics describe the evolution of the code and how developers interact with it. Temporal indicators, such as Age in Days [AGE], capture the maturity of a component. Churn-based measures quantify the extent of modifications performed over time, covering averages, maxima, and totals for both code churn and modified statements. Commit-based metrics capture the structure and magnitude of version control activity, including the total number of commits, commit-level line changes, added or deleted code, and repeated change patterns. Method-based measures capture fine-grained evolution through counts of method additions, deletions, and line-level modifications. Developer-based indicators quantify the social and organizational aspects of code evolution, such as the number of distinct authors, the distribution of contributions, and average developer experience, calculated as the average of all developers’ experience scores, where each score reflects

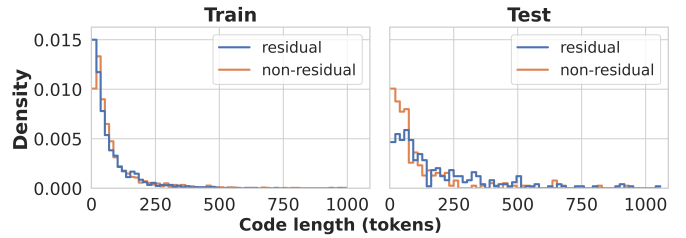


Fig. 2: Histograms of the train and test splits showcasing the distribution of token counts across residual and non-residual faults

how much, how long, and how recently a developer worked on the method. All product and process metrics were computed through purpose-built analysis scripts, validated across multiple runs to ensure consistency and correctness.

IV. EXPERIMENTAL SETUP

We conducted our experiments using both local development resources and large-scale computing infrastructure. Initial testing and model development were carried out on a standalone workstation equipped with a 12th Gen Intel Core i7-12700 processor (12 cores, 20 threads, up to 4.9 GHz), 32 GB of RAM, and an NVIDIA GeForce GTX 1660 SUPER GPU with 6 GB of dedicated memory. For the fine-tuning of deep learning models, we relied on NVIDIA A100 GPUs with 80 GB of memory; each training job was assigned 16 Intel Xeon CPU cores and 30 GB of RAM.

Building upon this computational setup, we prepared the data used for training and validation. We adopted a 90–10 random splitting strategy, allocating 90% of the samples to the training set and 10% to the validation set. The resulting dataset comprises 3570 faults in the training split (1770 residual and 1800 non-residual) and 450 in the validation split (213 residual and 217 non-residual). To evaluate generalizability beyond our primary dataset, we further employed BugsInPy [28], a benchmark of verified Python bugs. After classifying each bug, we obtained 284 residual and 221 non-residual faults. All splits were checked to ensure that no commit appears in more than one partition, preventing data leakage. In total, the final corpus contains 4525 faults, of which 2267 are residual and 2238 non-residual. The 90/10 split models a single-ecosystem scenario in which an organization trains the models on historical data and tests them on real-world, previously unforeseen use case scenarios. BugsInPy serves as a cross-project evaluation benchmark, as it contains entirely separate commits and code with no overlap with the training corpus. Specifically, the training corpus comprises 357 different projects, while BugsInPy contains 16 different projects. Out of these 16, 10 projects are also in the training corpus.

Fig. 2 provides an initial overview of the dataset by illustrating the distribution of code lengths, expressed in number of tokens, across residual and non-residual faults. Several trends emerge. First, code length in both training and test splits is right-skewed: while most samples are relatively short, a

TABLE II: Statement-level statistics by split and residual label.

Group	Unique Statements	Unique Tokens	Average Tokens
train (residual)	29,459	40,710	3.34
train (non-residual)	29,992	44,880	3.69
test (residual)	9,470	13,808	4.19
test (non-residual)	4,612	7,900	3.96

long tail of larger snippets appears in both classes. Second, in the training split, residual and non-residual distributions nearly overlap, with most examples below 250 tokens and only a minority exceeding 500. In contrast, the test split shows a slightly higher proportion of longer residual examples, reflecting the increased complexity of the faults in BugsInPy.

To complement these observations, Table II reports statement- and token-level statistics. In the training set, residual and non-residual faults contain a comparable number of unique statements (29k–30k), but non-residual examples include more unique tokens overall, leading to a higher average token count per statement. This distinction becomes even more pronounced in the test set, where residual cases display a higher average token density (4.19 vs. 3.96). Taken together, these findings indicate that residual bugs often appear in structurally dense code, while also confirming that the two classes remain well balanced in their structural properties.

V. EXPERIMENTAL RESULTS

In this section, we present the results of our experimental campaign on predicting *residual faults* in Python systems. Detecting residual faults is intrinsically challenging: these bugs often differ only subtly from pre-release defects, and no straightforward lexical or syntactic patterns clearly separate the two. Because of this, we investigate the task through a progressively more informed sequence of analyses, moving from general-purpose reasoning models to domain-adapted models and finally to metric-based learners.

To guide our evaluation, we define four research questions:

- **RQ1:** *Can code representations learned by LLMs discriminate residual from non-residual faults?*
- **RQ2:** *Can software metrics improve residual fault prediction?*
- **RQ3:** *Which software metrics contribute most to the classification?*
- **RQ4:** *Do software metrics and LLM-based code representations encode the same or complementary information?*

The order of our analyses follows a progression. We begin by assessing whether code representations learned by Large Language Models, including both proprietary systems and fine-tuned open-weight models, can capture patterns that distinguish residual from non-residual faults (RQ1). We then shift to an orthogonal perspective and examine whether traditional software metrics, such as those used in anomaly detection or supervised classification, provide a clearer signal for this prediction task (each model using the same fixed system

TABLE III: Performance of closed-source LLMs on the residual fault prediction task with 95% CI.

Model	Accuracy	Precision	Recall	F1
Gemini 3 Pro	0.574 _[.53,.62]	0.570 _[.53,.61]	0.986 _[.97,1.0]	0.723 _[.69,.76]
Claude 4.5 Sonnet	0.535 _[.49,.58]	0.615 _[.55,.68]	0.461 _[.40,.52]	0.527 _[.47,.58]
GPT-4.1	0.461 _[.42,.50]	0.531 _[.46,.60]	0.359 _[.30,.42]	0.429 _[.37,.48]

prompt, which instructed decisions of the supervised models, leveraging explainability techniques to characterize their contribution (RQ3). Finally, we investigate whether software metrics and code-based representations encode overlapping or complementary information, and thus whether combining them can improve residual-fault prediction (RQ4).

A. RQ1: Can code representations learned by LLMs discriminate residual from non-residual faults?

In this RQ, we investigate whether the internal representations of Large Language Models (LLMs) can distinguish between residual and non-residual faults. We divide our analysis into two parts: state-of-the-art closed-source models prompted in a zero-shot setting and open-weight models fine-tuned specifically for this task.

1) *Closed-Source General-Purpose LLMs:* We first evaluated three state-of-the-art proprietary models, namely Gemini 3 Pro, Claude 4.5 Sonnet, and GPT-4.1, on our test set as a baseline, using them through their respective APIs. Before running the evaluation, we verified that all three models understood the concept of a residual fault. When asked directly, each model correctly described it as a bug that escaped the testing phase, confirming that a lack of knowledge did not hinder the task itself. To ensure comparability, we evaluated each model using the same fixed system prompt, which defined a residual fault and asked them to classify the input code.

Their performance is summarized in Table III. The three models exhibit two main approaches that underscore the difficulty of the task and the absence of simple lexical or structural cues that these models can reliably exploit.

Gemini 3 Pro achieves the highest F1-score (0.72); however, its behavior is dominated by an extremely aggressive classification strategy. Out of the 284 residual faults in the test set, Gemini correctly classifies 280 as post-release (TP) and misses only 4 (FN), yielding a recall of 0.99 on this class. However, this comes at the cost of misclassifying almost every non-residual instance: among the 221 pre-release faults, it correctly flags only 10 as pre-release (TN) and incorrectly labels the remaining 211 as post-release (FP). This results in a misleadingly strong F1, driven not by discriminative ability but by systematically favoring the residual label. In practice, Gemini 3 Pro behaves as if “residual” were the default prediction.

Claude 4.5 Sonnet adopts a much more conservative stance. It produces predictions for 416 out of 505 test instances, abstaining on 89 cases. On the subset where it does emit a label, Claude correctly identifies 132 pre-release faults (TN) and produces 62 false positives, showing a noticeably better

ability than Gemini to recognize non-residual faults. At the same time, it correctly classifies only 84 residual faults (TP) and misclassifies 138 as pre-release (FN), yielding a recall of 0.38 on the residual class. This pattern indicates a reluctance to assign the residual label: Claude tends to default to pre-release unless the signal is overwhelmingly clear, which reduces false alarms but systematically misses many true residual faults.

GPT-4.1 falls in the same trap, remaining ineffective overall. Like Claude, its confusion matrix still shows a strong difficulty in capturing residual faults. Among the 221 pre-release cases, GPT-4.1 correctly identifies 131 (TN) and mislabels 90 as post-release (FP), achieving a precision comparable to Claude on the residual class. However, on the 284 residual cases it correctly classifies only 102 (TP) and misclassifies 182 as pre-release (FN), yielding a recall of 0.36 and an F1 of 0.43. Rather than favoring one class as strongly as Gemini, GPT-4.1 distributes errors across both, suggesting uncertainty. To assess whether the observed performance differences are statistically significant, we apply McNemar’s test [29], a non-parametric test designed to compare two classifiers evaluated on the same set of instances. The test operates on the 2×2 contingency table of discordant predictions, i.e., cases where one model is correct, and the other is not, and is the standard choice when per-sample predictions are available from a single test set [30]. Results show that GPT-4.1 differs significantly from both Gemini 3 Pro ($p=0.002$) and Claude 4.5 Sonnet ($p=0.019$), whereas Gemini and Claude do not differ significantly ($p=0.260$), indicating that, despite opposite classification biases, their overall error rates are comparable.

Closed Source LLMs on Code

Closed-source LLMs fail to predict residual faults reliably. Gemini 3 Pro inflates its F1 by labeling almost everything as residual, while Claude 4.5 Sonnet and GPT-4.1 rarely detect residual faults at all. None of the models learns a meaningful boundary, confirming that general-purpose LLMs are ineffective for this task.

2) *Fine-Tuned Deep Learning Models*: We next evaluate a set of trained deep learning models on the same prediction task. The goal here is twofold: (i) to assess whether fine-tuning on code can overcome the limitations observed with closed-source LLMs, and (ii) to determine whether code representations extracted by large pretrained models contain signals that generalize beyond project-specific context. For this purpose, we selected three representative architectures that cover the state-of-the-art DL models: *CodeT5+* for Classification [31] (770M parameters), an encoder–decoder model optimized for code understanding; *DeepSeek-Coder* [32] (1.3B parameters), a transformer decoder pretrained specifically on code repositories; and *CodeLlama* [33] (7B parameters), a decoder-only backbone specialized for code generation. Before training, each snippet was processed through our normalization pipeline, which removes comments and docstrings, anonymizes identifiers and literals, and cleans possible noise

TABLE IV: Performance of trained deep learning models on the residual fault prediction task with 95% CI.

Model	Accuracy	Precision	Recall	F1
CodeLlama	0.438 _[.40,.48]	0.500 _[.43,.57]	0.327 _[.27,.38]	0.396 _[.34,.45]
CodeT5+	0.406 _[.36,.45]	0.470 _[.41,.53]	0.440 _[.38,.50]	0.455 _[.40,.51]
DeepSeek-Coder	0.388 _[.35,.43]	0.426 _[.35,.50]	0.254 _[.20,.30]	0.318 _[.26,.37]

around the method (e.g., wrongly extracted keywords). This prevents models from relying on project-specific tokens, forcing them instead to use structural and semantic cues.

Despite their architectural differences, the three models exhibit limited predictive ability (Table IV). CodeT5+ achieves the highest F1 (0.46), followed by CodeLlama (0.40) and DeepSeek-Coder (0.32). However, none of the models reaches performance levels comparable to supervised metric-based approaches.

The confusion matrices reveal distinct behavioral profiles. CodeT5+ detects the most residual faults, correctly identifying 125 true positives, but at the cost of 141 false positives and still missing 159 residual cases (FN). CodeLlama trades recall for precision: it identifies 93 true positives with 93 false positives, yielding the highest accuracy (0.44) and precision (0.50) among the three, but missing 191 residual faults. DeepSeek-Coder performs worst overall, detecting only 72 residual faults while misclassifying 97 non-residual cases and leaving 212 residual faults undetected. McNemar’s test confirms that CodeLlama and CodeT5+ do not differ significantly ($p=0.202$), nor do CodeT5+ and DeepSeek-Coder ($p=0.460$), whereas CodeLlama significantly outperforms DeepSeek-Coder ($p < 0.001$).

Fine-tuned LLMs on Code

Fine-tuned deep learning models also struggle to predict residual faults. CodeT5+ achieves the best F1 (0.46) but at the cost of high false positives, while CodeLlama and DeepSeek-Coder miss the majority of residual cases. Despite architectural differences, none of the models extracts robust signals from code, confirming that residual fault prediction requires more specialized representations.

B. RQ2: Can software metrics improve residual fault prediction?

In this RQ, we investigate whether software metrics can help predict residual faults. To do so, we explicitly compare unsupervised anomaly detection algorithms with supervised classification models. All features were standardized to a zero mean and unit variance for both experimental settings.

1) *Unsupervised Anomaly Detection*: Unsupervised methods are appealing in this context because residual faults are relatively rare. For this setting, we use Isolation Forest (IF) [34], One-Class SVM (OC-SVM) [35], and Local Outlier Factor (LOF) [36]. They model expected behavior and flag statistically atypical cases, helping identify potential defects or quality anomalies [37]–[39].

TABLE V: Performance of unsupervised models with 95% CI.

Model	Accuracy	Precision	Recall	F1
IsolationForest	0.523 _[.48,.57]	0.642 _[.57,.72]	0.342 _[.29,.40]	0.446 _[.39,.50]
OneClassSVM	0.497 _[.45,.54]	0.650 _[.55,.74]	0.229 _[.18,.28]	0.339 _[.28,.40]
LocalOutlierFactor	0.493 _[.45,.54]	0.818 _[.70,.93]	0.127 _[.09,.17]	0.220 _[.16,.28]

TABLE VI: Performance of supervised models with 95% CI.

Model	Accuracy	Precision	Recall	F1
RandomForest	0.636 _[.59,.68]	0.622 _[.57,.67]	0.898 _[.86,.93]	0.735 _[.70,.77]
XGBoost	0.610 _[.57,.65]	0.611 _[.56,.66]	0.845 _[.80,.89]	0.709 _[.67,.75]
CatBoost	0.628 _[.59,.67]	0.622 _[.57,.67]	0.870 _[.83,.91]	0.724 _[.68,.76]

As shown in Table V, unsupervised anomaly detectors perform poorly on this task. Although they achieve reasonable precision on the residual class, they detect only a small portion of the true residual faults. IF correctly identifies 97 residual cases (TP) but misses 187 (FN), while OC-SVM and LOF degrade further to 64 and 35 true positives, respectively, leaving more than two hundred residual faults undetected (FN 220 and 249, respectively). In all three methods, false positives remain modest, but this stems from a consistent pattern: they overwhelmingly classify faults as non-residual, rarely flagging anomalies. McNemar’s test reveals no significant pairwise difference among the three detectors ($p \geq 0.13$ for all pairs), confirming that their poor performance is a structural limitation of the anomaly-detection paradigm for this task. These results confirm that residual faults do not manifest as outliers in metric space, and that unsupervised anomaly detectors are fundamentally unsuited for this prediction task.

Unsupervised ML Models on Software Metrics

Unsupervised anomaly detectors (IF, OC-SVM, LOF) fail to identify residual faults reliably. While they maintain reasonable precision, they suffer from extremely low recall, missing the vast majority of residual cases. Residual faults do not manifest as clear statistical outliers in the software metric space, rendering anomaly detection ineffective for this task.

2) *Supervised Classification*: Supervised models trained on the same metrics perform significantly better and reveal a different set of trade-offs (Table VI). RandomForest, XGBoost, and CatBoost reach F1-scores between 0.71 and 0.73 and show substantially improved sensitivity to residual faults. RandomForest, for example, correctly identifies 255 residual faults (TP: 255, FN: 29) but also produces 155 false positives (TN: 66). XGBoost and CatBoost behave similarly, obtaining 241–248 true positives while producing roughly 150 false positives each. Compared to the unsupervised baselines, these supervised learners recover more than three times as many residual faults and reduce false negatives, demonstrating that software metrics contain meaningful predictive signals. McNemar’s test shows no significant difference among the three supervised learners ($p \geq 0.11$ for all pairs), indicating that the

predictive signal resides in the metrics themselves rather than in the choice of classifier. In contrast, cross-family comparisons confirm that supervised models significantly outperform both unsupervised detectors (RandomForest vs. IsolationForest, $p < 0.001$) and closed-source LLMs (RandomForest vs. Gemini 3 Pro, $p = 0.003$), despite the latter’s nominally comparable F1.

When contrasted with the other deep learning models studied, the supervised ones occupy an interesting middle ground. They do not match Gemini 3 Pro’s extreme recall (which comes from labeling nearly everything as residual). Still, unlike Gemini, their high TP counts arise from actual discriminative patterns rather than a degenerate decision rule. At the same time, they outperform Claude 4.5 Sonnet and GPT-4.1, CodeLlama, CodeT5+, and DeepSeek-Coder in both TP and FN ratios, revealing a stronger ability to separate residual from non-residual cases than any model relying solely on representations learned from code. Yet their high false-positive rates reveal the limits of metric-only features.

Supervised ML Models on Software Metrics

Supervised learners trained on software metrics significantly outperform unsupervised detectors and many deep learning models. RandomForest, XGBoost, and CatBoost correctly identify 241–255 residual cases, reducing false negatives by an order of magnitude. However, all three models still generate several false positives: metric-based supervision captures useful predictive signals but struggles to achieve a clean separation between pre-release and post-release faults, positioning these models as safeguards rather than definitive tools.

C. *RQ3: Which software metrics contribute most to the classification?*

To identify the properties that distinguish residual from non-residual faults, we examine the feature-importance rankings produced by the three supervised learners from RQ2 and complement them with a SHAP-based directional analysis [40]. Fig. 3 reports the top-10 metrics for each model. The figure reports each model’s feature-importance scores, i.e., the model-specific estimates of how much each metric contributes to prediction, normalized so that importances sum to 1.

Across all models, it emerges that process metrics play a central role, in particular the age of the method (*AGE*), which is highly ranked in all three learners. The SHAP analysis indicates that higher *AGE* pushes the prediction towards the residual-fault class, suggesting that long-lived components are more likely to host faults that escape pre-release detection. In contrast, bug density (*BD*) and average commits per author (*ACA*) have an opposite effect: their SHAP means are negative, meaning that larger *BD* and *ACA* values are associated with a higher probability of non-residual faults. As expected, modules with many historical defects and intense patching activities have a reduced chance that further faults remain hidden.

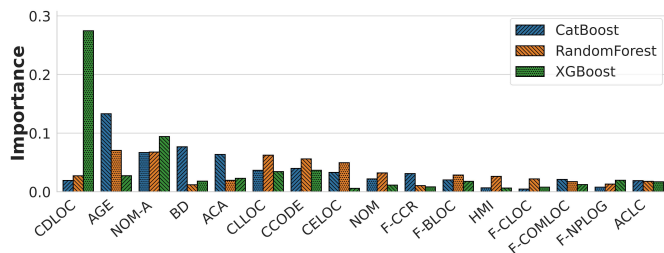


Fig. 3: Top-10 feature importances for the three supervised models used in RQ2

Structural and size-oriented class metrics are equally prominent. *CDLOC* (class declaration lines), *NOM-A* (number of declared methods), *CLLOC* (class LOC), *CCODE*, *CELOC*, and *NOM* all appear in the top-10 lists, especially for XGBoost and RandomForest. For these metrics the SHAP means are consistently positive, i.e., larger classes increase the likelihood of a residual-fault prediction.

File-level size and navigability measures provide additional explanatory power. *F-BLOC* (blank lines) and *F-NPLOG* (logarithm of paths) are important for XGBoost and, according to SHAP, higher values of both push towards the residual-fault class, indicating that large and structurally complex files are risk-prone. *F-COMLOC* (file comment lines) is relevant for CatBoost and also shows a positive SHAP mean, suggesting that highly documented components are not necessarily less risky. Conversely, *F-CCR* (comment-to-code ratio) tends to have a small but negative SHAP mean in CatBoost, showing that well-documented and written code is less prone to host residual faults.

While the three learners share these general tendencies, their emphases differ. XGBoost focuses on a small set of structural metrics (*CDLOC*, *NOM-A*, *CLLOC*, *F-BLOC*, *F-NPLOG*), all of which exhibit a clear positive SHAP direction. CatBoost distributes importance more evenly between process and structural dimensions, with *AGE*, *ACA* and *BD* among the most influential process metrics. RandomForest assigns similar relative importance to *AGE*, *NOM-A*, *CLLOC*, *CDLOC*, *CCODE*, and *CELOC*, but its SHAP means are close to zero, indicating weaker directional effects.

Key Software Metrics for Fault Prediction

Residual-fault prediction is driven by a combination of historical process indicators and specific structural size characteristics, with *AGE*, class/file size, and dependencies consistently associated with an increased residual-fault proneness. Post-release faults in Python mostly arise in older, larger, structurally complex, and badly documented components.

D. RQ4: Do software metrics and LLM-based code representations encode the same or complementary information?

Our results show that software metrics prove more useful than code representations for cross-project prediction tasks.

However, this does not imply that one representation is inherently superior. Our question is whether code can provide an additional, orthogonal source of information that can be combined with metrics to improve prediction when conditions allow. To explore this, we compared traditional software metrics with dense code representations extracted from the pre-trained CodeLlama model. Embeddings were obtained using its tokenizer and the penultimate transformer layer, followed by mean pooling to produce a single 4,096-dimensional vector per function. The model weights remained frozen, ensuring that the embeddings capture syntactic and semantic patterns learned during pretraining.

Using this representation, we examined the relationship between metrics and embeddings. After standardization, both spaces were reduced with PCA [41], yielding 29 components for metrics and 28 for embeddings, each preserving more than 95% of the variance. The association between the two is limited: the maximum absolute Spearman correlation between any metric and embedding component is 0.27, and the average maximum correlation per component is roughly 0.11–0.12. Canonical Correlation Analysis (CCA) [42] confirms this observation. Although the first few canonical correlations are moderate (about 0.71, 0.53, and 0.50), the mean across the first twenty is only around 0.24, indicating that the shared subspace is small. The geometric distinction between the two representations is evident in the Fig. 4, where components derived from metrics and embeddings form clearly separated clusters with almost orthogonal regression directions.

Taken together, these results suggest that improving testing does not require replacing metrics with embeddings or vice versa, but combining them. Metrics could help identify where additional testing and review should be directed, while embeddings capture semantic aspects of faults.

Differences between Code and Metrics

Metrics and embeddings share little information. After PCA, their components reach a maximum Spearman correlation of 0.27, typically around 0.11–0.12. CCA confirms this: the mean of the first twenty components is ≈ 0.24 . Their PCA clusters are nearly orthogonal, indicating that metrics capture structural factors while embeddings provide complementary semantic cues.

VI. DISCUSSION

Our results highlight a clear separation between code-based and metric-based approaches to residual fault prediction. Deep learning models such as CodeT5+, DeepSeek-Coder, and CodeLlama struggle to form a reliable decision boundary and consistently miss true residual faults. They also require substantial resources: training CodeT5+ and DeepSeek-Coder on an NVIDIA A100 with 80 GB of memory took about 1.3 hours, and CodeLlama took about 3 hours, with inference remaining slow and costly.

In contrast, traditional supervised models trained on software metrics are lightweight, fast, and more effective. On a

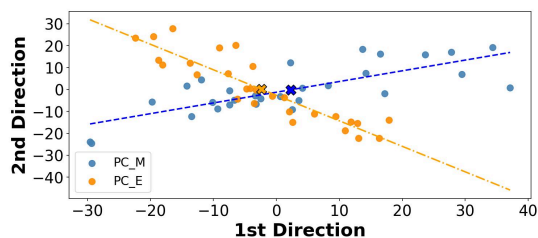


Fig. 4: PCA projection of metric and embedding components. The two X’s represent the centroids. The two clusters exhibit a near-orthogonality of the two representations.

standard server with a GTX 1660 Ti and an Intel Core i7-12700, they train in minutes and achieve higher recall. The precision of these models remains moderate, but prior work has shown that [43], [44] what matters in practice is the trade-off between inspection effort saved and residual faults caught. Our models retain 0.85–0.9 recall while substantially reducing the number of commits requiring inspection, effectively acting as safety filters that preserve most escaping defects at a manageable additional validation cost. Their interpretability further reveals the structural and historical patterns associated with fault survival, including age, size, churn, and complex modification histories. McNemar’s test shows no significant difference among RandomForest, XGBoost, and CatBoost ($p \geq 0.11$), indicating that the predictive power resides in the metrics rather than in a particular learning algorithm.

Metrics alone, however, do not capture all relevant information. Our representational analysis shows that metrics and embeddings encode largely independent signals: the maximum correlation between PCA components is only 0.27, typical values are around 0.11 to 0.12, and the mean canonical correlation across the first twenty components is about 0.24. The PCA projection forms two almost orthogonal clusters, indicating that embeddings carry semantic cues that metrics do not, even if those cues are insufficient on their own. These results point to a clear direction for future work. Models that combine structural indicators from metrics with semantic information from embeddings may achieve better stability and cross-project generalization. We conducted preliminary experiments with hybrid models that concatenate software metrics and code embeddings as input features. However, these models consistently produced degenerate behavior: when relying on metrics alone, performance remained stable, but whenever the models attempted to leverage the code representation, predictions degraded. We omitted detailed results for brevity, but the observed pattern suggests that the high dimensionality and variability of code embeddings overwhelm the compact metric signal during joint training, a challenge that warrants dedicated investigation with larger and more diverse datasets.

There remains a central question in our study: how can we improve the testing of residual faults? Our results suggest a two-part answer. First, testing should prioritize the components that metrics identify as high risk, such as old, large, highly modified, and frequently touched code, because

these structural and historical factors strongly correlate with fault survival. Second, testing should incorporate semantic cues, since metrics and embeddings capture complementary information. Combining structural indicators with semantic patterns offers a more complete view of where testing is most likely to miss defects that later manifest in production. Our statistical analysis supports this direction: the stability of supervised performance across three independent classifiers suggests that the metric signal is robust, while the orthogonality between metrics and embeddings ($\rho_{\max}=0.27$) indicates substantial room for improvement through hybrid representations. Modern LLMs are well-positioned to provide such semantic cues [45]–[47]: they have proved to hold great results when specifically fine-tuned on the same projects under test, but the challenge lies in capturing the extreme variability that the code representations have in cross-project scenarios.

VII. RELATED WORK

Software Fault Prediction (SFP) has long been a central research area in software engineering, aiming to identify defect-prone components early in the development lifecycle to guide testing and quality assurance. Foundational studies demonstrated the predictive power of static code metrics, focusing on size [48], control-flow complexity [49], and cognitive complexity [50]. These early works laid the groundwork for leveraging structural indicators to infer software quality.

As programming paradigms evolved, object-oriented metrics were introduced to capture higher-level design properties such as coupling, cohesion, and inheritance [51]. Empirical evidence confirmed the effectiveness of these metrics in predicting faults across a wide range of systems [52]. In parallel, process metrics, including code churn [53], change history [54], and developer activity [55], were adopted to reflect the dynamic aspects of software development. These metrics enhanced prediction models by incorporating historical and human-centric signals, often improving robustness [56].

Several systematic reviews have chronicled the evolution of SFP approaches. Catal et al. [57] emphasized the growing shift toward machine learning (ML)-based models and highlighted the need for publicly available datasets like NASA MDP [58] and PROMISE [59]. A follow-up study [60] further reported that method-level metrics and ML techniques were gaining traction due to their improved granularity and adaptability.

The emergence of DL has transformed the SFP landscape by enabling models to learn directly from raw source code. Akimova et al. [61] reviewed DL-based SFP methods, noting the growing relevance of attention mechanisms and transformer-based architectures. Qiao et al. [62] showed that DL models outperform traditional approaches by capturing syntactic and semantic patterns, while Feng et al. [63] demonstrated the utility of pretrained models like CodeBERT in learning long-range code dependencies. However, these models are opaque in their decision-making and demand significant computational resources, which may limit their adoption in production.

Alternative paradigms have also been explored to overcome the limitations of static models. Dynamic analysis techniques,

such as those presented by Zhang et al. [64], observe runtime behavior to detect execution-specific faults. Natural Language Processing (NLP) methods mine commit logs and bug reports for fault signals [65], and graph-based approaches model control-flow and data-flow dependencies using program graphs [66]. While effective in capturing rich structural and behavioral information, these methods often suffer from scalability issues and reduced transparency, making them less accessible for practical deployment in industry settings.

A persistent limitation in the SFP literature is the lack of focus on residual (post-release) faults. Most studies treat all defects as temporally undifferentiated, even when logs or timestamps are available [67], [68]. Cross-project studies [69]–[71] similarly assume a homogeneous fault distribution, ignoring the operational and contextual challenges associated with post-release failures. This modeling assumption restricts the applicability of such models in production, where residual faults can be the most costly and disruptive [6].

Recent work at Meta [45]–[47] demonstrates that risk models built on prior defects can effectively identify risky changes within a single organization, reducing outages. In that setting, models are trained and deployed on the same codebase with stable distributions of developers, processes, and dependencies. Our evaluation setting differs precisely in this regard: we train on one dataset and test on BugsInPy, an independent benchmark comprising different projects, developers, and histories. Under such cross-project distribution shifts, code-level representations become unreliable because structural idioms vary widely across repositories. Process and product metrics, by contrast, provide a compact and project-agnostic representation of development dynamics.

Our work targets residual fault prediction, explicitly focusing on predicting residual faults. It contrasts interpretable software metrics with code-based features and evaluates them in a cross-project setting. Our methodology emphasizes model transparency through metric-level interpretability, achieving strong predictive performance while providing new insights into the most effective features and models for post-release reliability assessment. To our knowledge, no prior work has isolated residual fault prediction as a distinct task and evaluated it across projects using multiple models.

VIII. THREATS TO VALIDITY

Language Scope. Our study focuses exclusively on Python systems, selected for their ubiquity in AI, scripting, and rapid-development ecosystems [72]–[74]. Python’s dynamic typing, and runtime-dependent behavior may increase the likelihood of defects surfacing only after deployment, potentially amplifying the prevalence and characteristics of residual faults. As a result, caution is required when generalizing our findings to statically typed or lower-level languages such as Java or C++, where type systems, compilation pipelines, and memory-safety constraints shape different fault distributions and testing regimes. Future work will extend our methodology to multi-language datasets to assess whether the patterns observed for Python residual faults hold across programming paradigms.

Classification Heuristic Bias Our residual vs. non-residual labeling relies on a heuristic that combines temporal evidence, issue metadata, and reporter identity. While this approach was deliberately designed to be conservative, favoring the unknown label when evidence is insufficient, it may introduce bias in the distribution of labels. In particular, commits without an associated issue are only classified as pre-release if they predate the first stable release, whereas no commit is labeled as post-release based solely on timestamp information. This choice avoids systematic misclassification due to non-linear release workflows, maintenance branches, or backports, but may underrepresent post-release faults in projects with sparse or inconsistent issue tracking. Although our manual validation demonstrates high accuracy ($\sim 95\%$) on a random sample, some residual noise may still be present in the automatically derived labels. We mitigate this risk by (i) publicly releasing the full labeling code and feature definitions, enabling inspection and replication, and (ii) conducting sensitivity analyses confirming that the main experimental trends persist even when commits labeled as unknown are excluded.

Model Coverage. We evaluated a representative selection of traditional ML models, fine-tuned deep learning models, and prompt-based LLMs. While these choices span the current spectrum of fault prediction techniques, each model has inherent architectural biases and differing capacities for code representation. Our goal was not to identify a universally superior model, but to investigate whether diverse paradigms can effectively leverage software metrics and code representations for residual fault prediction.

Software Metric Selection. Our metric set combines product, process, statistical, and Python-specific features drawn from prior defect-prediction literature. Excluding highly sparse or trivially correlated metrics may further influence model behavior. To confirm this, we conducted ablation analyses, confirming that the main performance trends persist after removing the redundant metrics [75]–[77].

IX. CONCLUSION AND FUTURE WORK

This paper presented an empirical study on residual fault prediction in Python software, combining software metrics with traditional ML and modern DL models. We constructed a balanced dataset of validated residual and non-residual faults to enable reliable evaluation. Our findings show that residual faults are primarily associated with historical and structural properties of the code. Process and coupling metrics consistently capture this signal. Supervised metric-based models achieve high recall but still produce false positives, making them practical safety filters for testing prioritization rather than definitive decision tools. Metrics and code embeddings capture complementary information, yet straightforward hybrid models do not consistently improve performance under cross-project settings, motivating more integration strategies.

Future work will focus on extending this approach to other programming languages and evaluating its applicability to large-scale, industrial software systems with complex development histories and proprietary constraints.

REFERENCES

- [1] H. Dolfin, "Project failure case study: Knight capital group," *Henrico Dolfin Blog*, June 2019. Accessed: 2025-11-21.
- [2] NYT, "How a self-driving uber killed a pedestrian in arizona," March 2018. Accessed: 2025-11-21.
- [3] I. S. Staff, "How the boeing 737 max disaster looks to a software developer," *IEEE Spectrum*, 2019. Accessed: 2024-11-21.
- [4] J. Herschmann, "Improve software quality by building digital immunity." <https://www.gartner.com/en/doc/735246-improve-software-quality-by-building-digital-immunity>, 2023. Accessed: 2025-04-20.
- [5] N. Nagappan, T. Ball, and A. Zeller, "Mining metrics to predict component failures," in *Proceedings of the 28th International Conference on Software Engineering, ICSE '06*, (New York, NY, USA), p. 452–461, Association for Computing Machinery, 2006.
- [6] R. Rwemalika, M. Kintis, M. Papadakis, Y. Le Traon, and P. Lorrach, "An industrial study on the differences between pre-release and post-release bugs," in *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pp. 92–102, 2019.
- [7] G. Giray, K. E. Bennis, Ö. Köksal, Ö. Babur, and B. Tekinerdogan, "On the use of deep learning in software defect prediction," *Journal of Systems and Software*, vol. 195, p. 111537, 2023.
- [8] A. Chowdhury, A. Hindle, and E. Shihab, "In situ defect prediction: Practical transfer learning for cross-project defect detection," *IEEE Transactions on Software Engineering*, 2024.
- [9] M. Cinque, D. Cotroneo, G. De Rosa, L. De Simone, and G. Farina, "Cosmos: A fault injection framework to assess hardware-assisted hypervisors," *IEEE Transactions on Dependable and Secure Computing*, 2025.
- [10] T. Hall, S. Beecham, D. Bowes, D. Gray, and S. Counsell, "A systematic literature review on fault prediction performance in software engineering," *IEEE Transactions on Software Engineering*, vol. 38, no. 6, pp. 1276–1304, 2011.
- [11] C. Catal and B. Diri, "A systematic review of software fault prediction studies," *Expert systems with applications*, vol. 36, no. 4, pp. 7346–7354, 2009.
- [12] TIOBE Software BV, "Tiobe index: The programming community index." <https://www.tiobe.com/tiobe-index/>. Accessed: 2025-11-27.
- [13] M. Rezaalipour and C. A. Furia, "An empirical study of fault localization in python programs," *Empirical Softw. Engg.*, vol. 29, June 2024.
- [14] R. Widayarsi, G. A. A. Prana, S. A. Haryono, S. Wang, and D. Lo, "Real world projects, real faults: evaluating spectrum based fault localization techniques on python projects," *Empirical Software Engineering*, vol. 27, no. 6, p. 147, 2022.
- [15] "Resource package." Zenodo, 2025. Zenodo record.
- [16] S. S. Rathore and S. Kumar, "A study on software fault prediction techniques," *Artificial Intelligence Review*, vol. 51, no. 2, pp. 255–327, 2019.
- [17] M. K. Thota, F. H. Shajin, and P. Rajesh, "Survey on software defect prediction techniques," *International Journal of Applied Science and Engineering*, vol. 17, no. 4, pp. 331–344, 2020.
- [18] GitHub, "Octoverse: AI leads Python to top language as the number of global developers surges." <https://github.blog/news-insights/octoverse/octoverse-2024/>, 2024. Accessed: Apr. 20, 2025.
- [19] Real Python, "Python News Roundup: November 2024." <https://realpython.com/python-news-november-2024/>, 2024. Accessed: Apr. 20, 2025.
- [20] S. Nanz and C. A. Furia, "A comparative study of programming languages in rosetta code," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, pp. 778–788, 2015.
- [21] H. Krasner, "The cost of poor software quality in the us: A 2022 report." <https://www.it-cisq.org/the-cost-of-poor-quality-software-in-the-us-a-2022-report/>, December 2022.
- [22] D. Cotroneo, G. De Rosa, and P. Liguori, "Pyresbugs: A dataset of residual python bugs for natural language-driven fault injection," in *2025 IEEE/ACM Second International Conference on AI Foundation Models and Software Engineering (Forge)*, pp. 146–150, IEEE, 2025.
- [23] S. Zafar, M. Z. Malik, and G. S. Walia, "Towards standardizing and improving classification of bug-fix commits," in *2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, pp. 1–6, 2019.
- [24] S. Levin and A. Yehudai, "Boosting automatic commit classification into maintenance activities by utilizing source code changes," in *Proceedings of the 13th International Conference on Predictive Models and Data Analytics in Software Engineering, PROMISE*, (New York, NY, USA), p. 97–106, Association for Computing Machinery, 2017.
- [25] M. Caulo and G. Scanniello, "A taxonomy of metrics for software fault prediction," in *2020 46th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, pp. 429–436, 2020.
- [26] SciTools, LLC, *Understand: The Software Developer's Multi-Tool*. SciTools, LLC, Lexington, KY, USA, 2025.
- [27] V. Raychev, P. Bielik, and M. Vechev, "Probabilistic model for code with decision trees," *SIGPLAN Not.*, vol. 51, p. 731–747, Oct. 2016.
- [28] R. Widayarsi, S. Q. Sim, C. Lok, H. Qi, J. Phan, Q. Tay, C. Tan, F. Wee, J. E. Tan, Y. Yieh, B. Goh, F. Thung, H. J. Kang, T. Hoang, D. Lo, and E. L. Ouh, "Bugsinpy: a database of existing bugs in python programs to enable controlled testing and debugging studies," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2020*, (New York, NY, USA), p. 1556–1560, Association for Computing Machinery, 2020.
- [29] Q. McNemar, "Note on the sampling error of the difference between correlated proportions or percentages," *Psychometrika*, vol. 12, no. 2, pp. 153–157, 1947.
- [30] T. G. Dietterich, "Approximate statistical tests for comparing supervised classification learning algorithms," *Neural Computation*, vol. 10, no. 7, pp. 1895–1923, 1998.
- [31] Y. Wang, L. Zhou, W. Chen, L. Dong, J. Wei, H. Wang, S. Li, D. Zhou, X. Xia, X. Lu, Q. Xie, H. Cheng, Z. Deng, X. Tan, P. Liang, Z. Zhang, Y. Li, X. Wang, H. Cheng, W. Cui, Y. Qin, X. Chen, E. Wong, Y. Lu, W. Yu, C. Pang, Z. Chen, Z. Xu, and W. Ye, "Codet5+: Open code large language models for code understanding and generation," *arXiv preprint arXiv:2305.07922*, 2023.
- [32] X. Wang, T. Sun, X. Yuan, Z. Feng, X. Hu, Z. Lin, Z. Zhu, Y. Ni, G. Shen, T. Xiao, and X. Xie, "Deepseek: Turning large language models into multi-turn semantic parsers," *arXiv preprint arXiv:2401.14196*, 2023.
- [33] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. C. Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve, "Code llama: Open foundation models for code," 2024.
- [34] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *Proceedings of the 2008 IEEE International Conference on Data Mining*, pp. 413–422, IEEE, 2008.
- [35] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural computation*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [36] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, "Lof: identifying density-based local outliers," in *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pp. 93–104, 2000.
- [37] Z. Ding and L. Xing, "Improved software defect prediction using pruned histogram-based isolation forest," *Reliability Engineering & System Safety*, vol. 204, p. 107170, 2020.
- [38] R. Moussa, D. Azar, and F. Sarro, "On the effectiveness of one-class support vector machine in different defect prediction scenarios," in *2024 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pp. 535–545, 2024.
- [39] R. Rabih, H. Vahdat-Nejad, W. Mansoor, and J. H. Joloudari, "Highly accurate anomaly based intrusion detection through integration of the local outlier factor and convolutional neural network," *Scientific Reports*, vol. 15, no. 1, p. 21147, 2025.
- [40] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in neural information processing systems*, vol. 30, 2017.
- [41] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley interdisciplinary reviews: computational statistics*, vol. 2, no. 4, pp. 433–459, 2010.
- [42] H. Hotelling, "Relations between two sets of variates," in *Breakthroughs in statistics: methodology and distribution*, pp. 162–190, Springer, 1992.
- [43] S. Herbold, "On the costs and profit of software defect prediction," *IEEE Transactions on Software Engineering*, vol. 47, no. 11, pp. 2617–2631, 2019.

- [44] S. Tunkel and S. Herbold, "Exploring the relationship between performance metrics and cost saving potential of defect prediction models," *Empirical Software Engineering*, vol. 27, no. 7, p. 182, 2022.
- [45] A. Mockus, P. C. Rigby, R. Abreu, A. Akkerman, Y. Bhootada, P. Bhuptani, G. Ghardhora, L. H. Dao, C. Hawley, R. He, *et al.*, "Code improvement practices at meta," *arXiv preprint arXiv:2504.12517*, 2025.
- [46] A. Mockus, R. Abreu, P. C. Rigby, D. Amsellem, P. Bansal, K. Chinniah, B. Ellis, P. Fan, J. Ge, B. He, *et al.*, "Leveraging risk models to improve productivity for effective code un-freeze at scale," *ACM Transactions on Software Engineering and Methodology*, vol. 34, no. 7, pp. 1–24, 2025.
- [47] R. Abreu, V. Murali, P. C. Rigby, C. Maddila, W. Sun, J. Ge, K. Chinniah, A. Mockus, M. Mehta, and N. Nagappan, "Moving faster and reducing risk: Using llms in release deployment," in *2025 IEEE/ACM 47th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*, pp. 448–457, IEEE, 2025.
- [48] F. Akiyama, "An example of software system debugging," *Proceedings of the IFIP Congress*, pp. 353–359, 1971.
- [49] T. J. McCabe, "A complexity measure," *IEEE Transactions on Software Engineering*, vol. SE-2(4), pp. 308–320, 1976.
- [50] M. H. Halstead, "Elements of software science," *Elsevier*, 1977.
- [51] S. R. Chidamber and C. F. Kemerer, "A metrics suite for object-oriented design," *IEEE Transactions on Software Engineering*, vol. 20(6), pp. 476–493, 1994.
- [52] L. C. Briand, J. Wüst, J. W. Daly, and D. V. Porter, "Exploring the relationships between design measures and software quality in object-oriented systems," *Journal of Systems and Software*, vol. 51(3), pp. 245–273, 2000.
- [53] T. L. Graves, A. F. Karr, J. S. Marron, and H. Siy, "Predicting fault incidence using software change history," *IEEE Transactions on Software Engineering*, vol. 26(7), pp. 653–661, 2000.
- [54] N. Ohlsson and H. Alberg, "Predicting fault-prone software modules in telephone switches," *IEEE Transactions on Software Engineering*, vol. 22(12), pp. 886–894, 1996.
- [55] A. E. Hassan, "Predicting faults using the complexity of code changes," *Proceedings of ICSE*, pp. 78–88, 2009.
- [56] D. Radjenović, M. Heričko, R. Torkar, and A. Živković, "Software fault prediction metrics: A systematic literature review," *Information and Software Technology*, vol. 55(8), pp. 1397–1418, 2013.
- [57] C. Catal and B. Diri, "A systematic review of software fault prediction studies," *Expert Systems with Applications*, vol. 36(4), pp. 7346–7354, 2009.
- [58] N. M. D. Program, "Nasa metrics data program (mdp) datasets." <http://mdp.ivv.nasa.gov/>. Accessed: 2025-04-30.
- [59] G. Boetticher, T. Menzies, and T. Ostrand, "Promise repository of empirical software engineering data." <https://github.com/ApoorvaKrisna/NASA-promise-dataset-repository>, 2005. Accessed: 2025-04-30.
- [60] C. Catal, "Software fault prediction: A literature review and current trends," *Expert Systems with Applications*, vol. 38(4), pp. 4626–4636, 2011.
- [61] E. N. Akimova, A. Y. Bersenev, A. A. Deikov, K. S. Kobylkin, A. V. Konygin, I. P. Mezentsev, and V. E. Misilov, "A survey on software defect prediction using deep learning," *Mathematics*, vol. 9(11), p. 1180, 2021.
- [62] L. Qiao, X. Li, Q. Umer, and P. Guo, "Deep learning based software defect prediction," *Neurocomputing*, vol. 385, pp. 100–110, 2020.
- [63] Y. Feng, S. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, and D. Jiang, "Codebert: A pre-trained model for programming and natural languages," *Findings of EMNLP*, pp. 1536–1547, 2020.
- [64] F. Zhang, S. Kim, and S. Thummalapenta, "Faulttracer: A fault localization-based approach to fixing bugs," in *2013 20th Working Conference on Reverse Engineering (WCORE)*, pp. 272–281, 2013.
- [65] H. Hata, O. Mizuno, and T. Kikuno, "Bug prediction based on fine-grained module histories," in *Proceedings of the 34th International Conference on Software Engineering (ICSE)*, pp. 200–210, 2012.
- [66] Z. Li, X. Mao, and L. Zhang, "Deep learning-based bug detection in source code using graph embedding," in *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 143–153, 2017.
- [67] S. Wang, T. Liu, J. Nam, and L. Tan, "Deep semantic feature learning for software defect prediction," *IEEE Transactions on Software Engineering*, vol. 46, no. 12, pp. 1267–1293, 2020.
- [68] S. D. Palma, D. D. Nucci, F. Palomba, and D. A. Tamburri, "Within-project defect prediction of infrastructure-as-code using product and process metrics," *IEEE Transactions on Software Engineering*, vol. 48, no. 6, pp. 2086–2104, 2022.
- [69] Z. Li, Y. Jin, and H. He, "Cross-project defect prediction via transfer learning: A benchmark study," *Empirical Software Engineering*, vol. 26(4), pp. 1–37, 2021.
- [70] Z. Liu, T. Su, M. A. Zakharov, G. Wei, and S. Lee, "Software defect prediction based on residual/shuffle network optimized by upgraded fish migration optimization algorithm," *Scientific Reports*, vol. 15, p. Article 7201, 2025.
- [71] F. Lomio, S. Moreschini, and V. Lenarduzzi, "A machine and deep learning analysis among sonarqube rules, product, and process metrics for faults prediction," *Empirical Software Engineering*, vol. 27, no. 3, p. 189, 2022.
- [72] B. Ray, D. Posnett, V. Filkov, and P. Devanbu, "A large-scale study of programming languages and code quality in github," in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, 2014.
- [73] M. Beller, G. Gousios, and A. Zaidman, "Oops, my tests broke the build: An explorative study of travis ci with github," in *Proceedings of the 14th International Conference on Mining Software Repositories*, 2016.
- [74] M. Pradel and K. Sen, "Deepbugs: A learning approach to name-based bug detection," in *Proceedings of the ACM on Programming Languages*, 2018.
- [75] X. Yang and W. Wen, "Ridge and lasso regression models for cross-version defect prediction," *IEEE Transactions on Reliability*, vol. 67, no. 3, pp. 885–896, 2018.
- [76] J. Jiarpakdee, C. Tantithamthavorn, A. Ihara, and K. Matsumoto, "A study of redundant metrics in defect prediction datasets," in *2016 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)*, pp. 51–52, 2016.
- [77] J. Jiarpakdee, C. Tantithamthavorn, and C. Treude, "The impact of automated feature selection techniques on the interpretation of defect models," *Empirical Softw. Engg.*, vol. 25, p. 3590–3638, Sept. 2020.