

Improved IR-based Bug Localization with Intelligent Relevance Feedback

Asif Mohammed Samir
Department of Computer Science
Dalhousie University
Halifax, NS, Canada
asifsamir@dal.ca

Mohammad Masudur Rahman
Department of Computer Science
Dalhousie University
Halifax, NS, Canada
masud.rahman@dal.ca

Abstract—Software bugs pose a significant challenge during development and maintenance, and practitioners spend nearly 50% of their time dealing with bugs. Many existing techniques adopt Information Retrieval (IR) to localize a reported bug using textual and semantic relevance between bug reports and source code. However, they often struggle to bridge a critical gap between bug reports and code that requires in-depth contextual understanding, which goes beyond textual or semantic relevance. In this paper, we present a novel technique for bug localization—*BRaIn*—that addresses the contextual gaps by assessing the relevance between bug reports and code with Large Language Models (LLM). It then leverages the LLM’s feedback (a.k.a., Intelligent Relevance Feedback) to reformulate queries and re-rank source documents, improving bug localization. We evaluate *BRaIn* using a benchmark dataset—*Bench4BL*—and three performance metrics and compare it against six baseline techniques from the literature. Our experimental results show that *BRaIn* outperforms baselines by 87.6%, 89.5%, and 48.8% margins in MAP, MRR, and HIT@K, respectively. Additionally, it can localize $\approx 52\%$ of bugs that cannot be localized by the baseline techniques due to the poor quality of corresponding bug reports. By addressing the contextual gaps and introducing Intelligent Relevance Feedback, *BRaIn* advances not only theory but also improves the IR-based bug localization.

Index Terms—Bug Localization, Query Reformulation, Intelligent Relevance Feedback, Information Retrieval, Large Language Models, Natural Language Processing, Software Engineering

I. INTRODUCTION

Software bugs can cause major financial losses and lead to data breaches, security vulnerabilities, and operational disruptions [1], [2]. A recent software bug from Microsoft-owned CrowdStrike caused several hours of disruption in the U.S. airline industry, nearly halting operations and resulting in over \$10 billion in damages [3], [4]. Developers at the major IT companies, such as Microsoft and Google, have reported bug resolution as a top concern [5]. According to existing studies up to 50% of the programming time is spent by developers on finding, understanding, and fixing software issues [6]–[8]. Thus, any automated support to tackle these challenges can greatly benefit the developers.

Software bugs are submitted to bug-tracking systems (e.g., Bugzilla, JIRA) as bug reports, which might capture crucial hints for resolving software-related issues. Developers often rely on these reports to trace the origin of bugs in the code. However, the content and quality of bug reports can vary

significantly based on their submitters’ level of expertise and articulation skills. In particular, there might be variations in word choice and presence of technical terms [9], [10]. Such variations pose challenges for developers when pinpointing the root cause of defects, even for seasoned practitioners [11]. To address these challenges, there has been significant research targeting the detection or localization of software bugs over the last few decades.

Researchers have presented two major categories of methods to automatically localize software bugs: program spectrum analysis and Information Retrieval. First, spectrum-based methods rely on program execution traces for fault localization. However, the execution traces are not always readily accessible, which makes these methods less scalable [12], [13]. On the other hand, Information Retrieval (IR)-based methods use overlapping terms or keywords between bug reports and source code to localize bugs [14]–[18]. They are lightweight and scalable. However, they also struggle with the *vocabulary mismatch problems* [9] and may not always deliver satisfactory results due to sporadic term matching. Researchers have also incorporated historical data from past bug reports, code change history, past bug fixes, and bug recurrences [19], [20]. Although these enhancements have been reported to improve the performance of the IR-based methods in localizing bugs, a recent study [21] suggests that they do not significantly outperform the previous methods.

Recent IR-based techniques focus on search queries and attempt to improve their queries by capturing syntactic, co-occurrence, and hierarchical dependencies among the words in bug reports [10], [11], [22], [23]. However, these methods only use terms found in bug reports, which could be poorly written or insufficient [11]. As a result, they frequently fail to bridge the gap between natural language from bug reports and programming code from a project when searching for software bugs. To address this issue, several techniques attempt to enhance queries with relevant terms extracted from source documents through relevance feedback mechanisms [10], [24]–[29]. However, the majority of these techniques naively consider the top few documents (based on textual similarity) as relevant, overlooking the need for a comprehensive understanding of the code. As a result, they may not always capture the most meaningful terms from source code

for their search queries. [22], [27]. Thus, the existing IR-based techniques for bug localization suffer from two major challenges as follows.

(a) Relevance feedback against search queries might not be always relevant: Gay et al. [29] proposed a manual, iterative approach that leverages relevance feedback from developers and constructs queries to search for buggy source documents. In contrast, Sisman et al. [30] and Kim et al. [27] select the top few documents as relevant (a.k.a., pseudo relevance feedback) and leverage the feedback to improve their search queries. However, these techniques rely heavily on textually similar documents, which may not be always relevant, especially when dealing with source code and bug reports. Thus, a deeper understanding of both bug reports and source code is warranted to improve the relevance feedback mechanism and the subsequent steps of Information Retrieval (e.g., query reformulation, retrieval).

(b) Textual and semantic relevance might not be sufficient: Bug reports contain not only natural language texts but also technical jargons, commit diffs, stack traces, and program elements [10]. These artifacts describe the context and symptoms of encountered bugs [22]. Since natural language is loosely structured, it can introduce ambiguity by expressing the same idea in various ways [9]. Similarly, programming languages are more structured yet allow syntactically diverse expressions (e.g., iterative vs. functional approaches) and arbitrary naming conventions [30]–[32]. This flexibility can result in textual mismatches, where keywords or key phrases in the bug report (e.g., “*download failed*”) do not directly match the identifiers in the code (e.g., `fetchResource`). At the same time, semantic mismatches can arise when a problem encountered in the bug report does not correspond to the programming task implemented in the code. For example, the encountered problem – “*download failed*” – might not align well with the task – “*HTTP/FTP operation task and get packets*” if the word-level semantics are considered only. It requires an understanding of the relationship between network operations and file downloading to establish their connection. In other words, to localize such bugs, automated tools or methods need to go beyond surface-level matching and comprehensively understand the context of an encountered problem as well as the functionality of the corresponding source code.

In this paper, we present a novel technique – *BRaIn* – to support bug localization using Information Retrieval (IR) and Intelligent Relevance Feedback (IRF). Our approach overcomes the challenge of contextual understanding of software bugs using Transformer models [33] and localizes the bugs leveraging such understanding. First, *BRaIn* collects potentially buggy documents from a codebase by analyzing their contextual relevance to a bug leveraging transformer models (e.g., Mistral [34]). That is, unlike the existing methods, our method captures more human-like feedback to a query (a.k.a., Intelligent Relevance Feedback). Second, it extracts appropriate terms from these documents and expands the original query by further leveraging the captured feedback.

Finally, *BRaIn* reranks the source documents by executing the expanded query and employing the relevance feedback, providing a refined list of suspicious source documents.

We conducted experiments using 4,683 bug reports from a benchmark dataset– Bench4BL [21]. We evaluated the performance of our approach using three commonly used metrics: Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), and HIT@K. Our approach is compared with six suitable baselines from the literature [10], [30], [35]–[38]. *BRaIn* consistently outperformed existing techniques, showing 19.3% and 87.6% higher MAP scores than that of traditional and Machine Learning (ML)-based approaches, respectively. Similar gains were observed in MRR (17.5% and 89.5%) and HIT@10 (12.2% and 48.8%). These results underscore the effectiveness and superiority of our proposed technique in software bug localization.

Thus, this research make following contributions-

- A novel relevance feedback mechanism, namely *Intelligent Relevance Feedback (IRF)*, that leverages the code understanding and reasoning capability of the LLM to offer useful feedback to a search query. It is neither naive like pseudo-relevance feedback nor costly like human feedback.
- A novel approach – *BRaIn* – that localizes software bugs using effective search queries and retrieval, supported by the intelligent relevance feedback mechanism.
- An extensive evaluation of *BRaIn* using three commonly used metrics and $\approx 4.7K$ bug reports and comparison with six baselines from three areas of the literature.
- A replication package* with a prototype, a curated dataset, and configuration details for third-party replication and reuse.

II. MOTIVATIONAL EXAMPLE

In this section, we present a motivating example to demonstrate the benefits of our proposed technique for bug localization. Let us consider the example bug report in Table I that discusses access problems to an LDAP server. The bug manifests as a failure in the authentication process, where the system returns an HTTP code of 403 (a.k.a., forbidden) instead of prompting for necessary credentials. This behavior results in a denial of access to the LDAP services and hinders any migration to a newer version of the services.

Fig. 1 presents the source code triggering the bug. The root cause of this bug is a subtle omission in the code handling authentication process. We see that the switch statement in the buggy version of code fails to account for the BASIC authentication type. Instead, it handles the PLAIN authentication type, which is semantically closer but not equivalent. On the other hand, the bug report mentions “BASIC” HTTP authentication, which is not present in the target code. This terminology mismatch creates a disconnect between the high-level system behavior described in the bug report and the code level implementation, making the detection of bugs challenging. As a result, traditional text-based search methods

TABLE I: An Example of Bug Report and Search Techniques

Bug ID# 2013 (Wildfly CORE)		Rank
Title	Unable to access HTTP management interface secured by legacy LDAP realm.	72
Description	When the HTTP management interface is secured with a legacy security realm using LDAP, the user is not prompted to provide credentials as should be in the case of BASIC HTTP authentication mechanism. Instead, a 403 HTTP status is returned directly. Users won't be able to migrate their current (6.4, 7.0) configuration to 7.1 without change.	13
Baseline Query	Bug Title + Bug Description	24
NextBug [37]	Cosine Similarity (Embedding + TF-IDF) between Bug Report and Source Documents	25
<i>BRaIn</i>	Intelligent Relevance Feedback + Query Expansion + Scoring	1

```

@@ -236,7 +236,7 @@ public class SecurityRealmService
private AuthMechanism toAuthMechanism
    (String mechanismType, String mechanismName) {
        switch (mechanismType) {
            case "SASL": ...
            case "HTTP":
                switch (mechanismName) {
                    case "DIGEST":
                        return AuthMechanism.DIGEST;
-                   case "PLAIN":
+                   case "BASIC":
                        return AuthMechanism.PLAIN;
                }
            break;
        }
        return null;
    }
}

```

Fig. 1: Buggy Code with Diff

perform poorly and retrieve the buggy code at 72nd, 13th, and 24th positions when title, description, or their combination are used as queries, respectively (Table I). Even after employing embedding-based semantic relevance, NextBug [37] struggles to link the code to the bug, placing it at 25th position.

The above evidence suggests that an in-depth analysis involving contextual relevance is essential. A seasoned developer would recognize the missing clause of BASIC HTTP authentication in the code, although it is not explicitly stated in the bug report. By probing deeper, they would infer that the missing BASIC authentication is likely the root cause of the reported issue.

Large language models (e.g., Mistral) are exposed to a vast amount of data, including text and code. As a result, they can identify patterns (as humans do) and infer missing details, making them adept at handling tasks that require deep contextual understanding. As shown in Table I, BRaIn leverages such capabilities to obtain intelligent feedback against its query, reformulates the query, and returns the buggy code as the topmost result by executing the query against an IR-based method (e.g., BM25 [39]).

III. METHODOLOGY

Fig. 2 shows the schematic diagram of our proposed technique – BRaIn – for software bug localization. We discuss its different steps in the following section.

A. Document Indexing and Retrieval

1) *Indexing*: To detect software bugs using Information Retrieval (IR), the first step is to index the source code documents from a code repository. We chose Elasticsearch [40] for indexing due to its reliability, support for diverse data types, and easy integration with computing systems (e.g., cloud). We collected 45 subject systems from an existing benchmark dataset – Bench4BL [21] – and indexed the source code (Step 1, Fig. 2) from 684 buggy versions of these systems. Our idea was to detect a bug in the exact version of the software system stated in the corresponding bug report. During the indexing, we employed Elasticsearch’s default analyzer to perform common pre-processing operations (e.g., tokenization, lowercase conversion, and removal of stop words).

2) *Retrieval of Potentially Buggy Documents using Textual Relevance*: To retrieve potentially buggy documents, we use bug reports (i.e., bug title and description) as queries (Step 2, Fig. 2). When we pass these queries to Elasticsearch, it preprocesses them using the standard analyzer and returns the top-K (e.g., 50) results through query execution. To narrow down the search results, we also apply additional filters, such as system and version information from each bug report. Without these filters, the retrieved documents could be irrelevant or noisy. This step provides a set of source documents ranked by their textual relevance against a bug report by employing Elasticsearch’s default retrieval algorithm, Okapi BM25 [39].

B. Intelligent Relevance Feedback

Once we have the results from Elasticsearch, we employ advanced prompt methods and Large Language Models (LLM) to determine the relevance between a bug report and each result (i.e., source code). LLMs have shown remarkable capabilities understanding natural language texts and source code [41]–[43]. We leverage their capabilities to capture intelligent relevance feedback against a query (a.k.a., bug report) as follows. To achieve this, we use prompt engineering, document segmentation, and finally relevance estimation as follows.

Prompt Engineering: *Prompting* is a novel method that instructs the LLMs (e.g., Mistral) to generate meaningful responses without any expensive training [44]–[46]. It involves crafting appropriate instructions to guide LLM outputs and make them applicable to different problem-solving tasks [47]–[51]. LLMs have been found to be effective with well-designed prompts that are clear, specific, and actionable [45], [52]. We first developed a candidate prompt to determine whether a given code segment triggers a reported bug. It instructs the LLM to find the relevance, deliver the output in a JSON format, and act as a rational software engineer, incorporating the contextual information from the bug report and code segment.

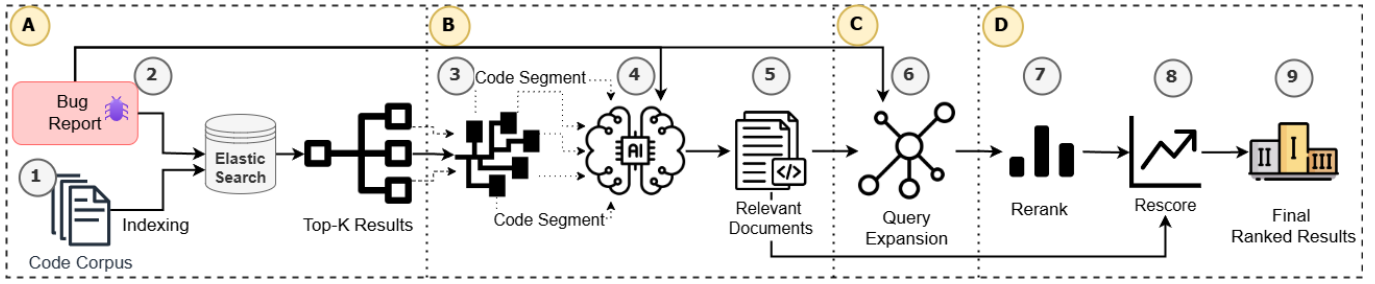


Fig. 2: Schematic Diagram of *BRaIn*:

(A) Document Indexing & Retrieval, (B) Intelligent Relevance Feedback, (C) Query Expansion, and (D) Bug Localization

To refine our candidate prompt, we employed SAMMO [53], a framework that explores and evaluates various prompt configurations using beam search. We configured SAMMO with LLaMA-3 [54] and used a small dataset of 20 bug reports with corresponding buggy code segments (ground truth) to guide our optimization process. SAMMO iteratively generated prompt variants by applying various modification operations to the candidate prompt with LLaMA’s assistance. In each iteration, we used LLaMA, bug reports and ground truth code to determine the fitness of each prompt and provide a performance update to SAMMO. Through an extensive search, SAMMO was able to find an optimized version of the prompt. Table II shows the optimized prompt template, used in the subsequent steps of our technique. Both candidate and optimized prompts can be found in the replication package.

Segmentation: We divide the source code documents from Elasticsearch into smaller segments to determine their relevance using prompting and LLM (Step 3, Fig. 2). According to an existing work [55], breaking up texts into smaller segments helps the attention mechanism focus on specific parts, which could be useful for our relevance estimation task. In our work, we adopt a simple method to capture code segments rather than collecting program slices. The slicing methods often create slices that are either too small to capture meaningful context or too large, introducing irrelevant contexts and potentially exceeding the token limits of the LLMs [56]. Therefore, we used a widely adopted library for static analysis – JavaParser [57], to extract code segments such as methods, constructors, interfaces, and enums from a document.

Determining the relevance of code: To determine code relevance, we employ LLaMA [54], Mistral [34], and Qwen [58] models in a zero-shot setting, and provide a bug report and a code segment (e.g., method, constructor) as *context* (Step 4, Fig. 2), respecting the token limits of these models (e.g., 8,192). We used the optimized template in Table II for LLaMA and Qwen, consisting of three key elements: *system*, *user*, *assistant*. On the other hand, the *system* element does not apply to Mistral, and thus its prompt template was adapted accordingly. We use Hugging Face’s [59] AutoTokenizer to perform model-specific formatting of the prompt and vLLM [60] to parallelize computations across the GPU in batches, increasing throughput. Each employed model against our context

provides a response. For example, for the showcase bug report (Table I) and corresponding buggy code (Fig. 1), we obtained the JSON response `{"relevance": "yes"}` from the Mistral. We also found a small number of cases where the outputs are malformed or incomplete JSON. In such instances, we perform string matching within the response (e.g., yes, no) to capture the relevance estimate of the code by the LLM. The relevance estimates of the code segments (collected from the results of Elasticsearch) serve as an intelligent feedback by the LLM to the original query. We coin this as *Intelligent Relevance Feedback (IRF)*.

TABLE II: Prompt Template for Relevance Feedback

<p>System: You are a helpful AI software engineer specializing in identifying buggy code segments given a bug report. Analyze the provided bug report and the JAVA code segment to determine if the code segment is responsible for causing the bug described in the bug report. You need to understand the functionality of the code segment and the details of the bug report to determine the relevance of the code segment to the bug report. There are two possible outputs: ‘yes’, ‘no’. - ‘yes’: The code is responsible for the bug described in the bug report. - ‘no’: The code is NOT responsible for the bug described in the bug report. Provide your output in JSON format like this sample: <code>{"relevance": "yes"}</code>. Act like a rational software engineer and provide output. Avoid emotion and extra text other than JSON.</p> <p>User: Analyze the following bug report and code segment: Bug Report: <BUG REPORT> Code Segment: <CODE SEGMENT></p> <p>Please determine if the code segment is responsible for the bug described in the bug report.</p> <p>Assistant:</p>

C. Query Expansion

Using the Intelligent Relevance Feedback (IRF), we expand an original query (Step 6, Fig. 2). Unlike the earlier work that relies on pseudo-relevance feedback [27], [29], we choose the source code documents marked as relevant by the LLM for query expansion. In pseudo-relevance feedback, the top few documents retrieved by an IR method are naively considered as relevant and are used for query expansion. On the other hand, our underlying idea is that documents contextually relevant to

the bug reports (i.e., IRF) provide terms that can complement the original query. We adapt an existing work of Rahman et al. [61] to capture appropriate terms from the relevant source documents as follows.

First, we parse each of the source documents retrieved by Elasticsearch and extract class, method, and field signatures from them. These signatures capture the intent of the code, whereas the detailed implementation code could be noisy [61]. We extract the signatures using a lightweight Python library – Javalang [62]. Next, we split camel case tokens from these signatures and turn them into textual phrases by combining their split tokens. We preprocess each of the phrases by filtering out stop words and programming keywords. Then, we construct a term graph $G(V, E)$ by encoding terms as vertices (V) and co-occurring terms as connecting edges (E). Subsequently, we apply the PageRank algorithm [63] in Eq. 1 to the term graph to select the influential terms.

$$PR(V_i) = \frac{1-d}{N} + d \sum_{V_j \in M(V_i)} \frac{PR(V_j)}{L(V_j)} \quad (1)$$

Here, $PR(V_i)$ represents the PageRank of vertex V_i . The term d denotes the damping factor with a default set to 0.85. N refers to the total number of vertices in the graph, while $M(V_i)$ indicates the set of vertices linked to V_i . Finally, $PR(V_j)$ and $L(V_j)$ represent the PageRank and outbound links of vertex V_j .

The algorithm assigns an initial score to each vertex (V_i) and iteratively updates it, prioritizing the vertices with higher connections. This process repeats until the scores stabilize or the algorithm reaches its maximum iteration limit (e.g., 100). Once the computation is done, we select the top- N (e.g., 10) weighted terms returned by the algorithm [61]. Finally, we expand the original query (i.e., bug report) with these terms, complementing bug reports with contextually relevant terms from source code, leveraging intelligent relevance feedback.

D. Bug Localization

We leverage our expanded query above and Intelligent Relevance Feedback (IRF) from the LLM to localize the buggy source documents as follows.

1) *Reranking*: We determine the relevance between each result from Elasticsearch and our expanded query employing BM25 algorithm, and rerank them according to their relevance (Step 7, Fig. 2). This step provides us with ranked results and their BM25 scores. The updated ranks could be useful since the expanded query contains more meaningful terms from the source documents.

2) *Rescoring*: We also enhance the ranking of source documents by incorporating IRF from LLM into their scores (Step 8, Fig. 2). First, we normalize the BM25 scores of the retrieved documents using a softmax function [64], producing a set of scores that add up to 1. The softmax function amplifies the differences among its input values exponentially, making their difference clearer. Given that the BM25 scores of the documents could have high variance, this allows the softmax function to highlight the textually relevant results.

However, since textual relevance might not be sufficient, we also leverage the LLM’s feedback against each result document. To incorporate this, we promote the relevant documents and penalize the irrelevant ones, marked by the LLM. This combined approach (Eq. 2) incorporates both textual and contextual relevance in the document ranking as follows.

$$score_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}} \cdot r_i, \quad r_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ doc is relevant} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Here, $score_i$ represents the score of the i^{th} document. The term r_i denotes the binary relevance feedback, indicating whether the document is relevant or not (details in Section III-B). The terms z_i and z_j in the softmax function correspond to the BM25 scores of documents i and j .

Finally, we rank the source documents based on their scores for their potential to be buggy (Step 9, Fig. 2) and return the top- K (e.g., $K=10$) documents. Our scoring process aims to bridge the gap between bug reports and source code by incorporating a deeper contextual understanding of the LLM and going beyond their textual and semantic relevance (Table I).

TABLE III: Dataset

(a) Dataset Summary			(b) Train-Test Split		
Project	Systems	Bug Reports	Project	Train	Test
Spring	25	1,802	Spring	1,429	373
Apache	25	1,802	Apache	1,246	313
Wildfly	5	806	Wildfly	643	163
Commons	8	507	Commons	402	105
JBoss	1	9	JBoss	7	2
Total	42	4,683	Total	3,727	956

IV. EXPERIMENTS

We curate a dataset of $\approx 4.7K$ bug reports from the benchmark dataset Bench4BL and evaluate using three appropriate metrics from the relevant literature — Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), and HIT@K ($K=1, 5, 10$) [10], [35]. We experiment with three different LLMs and compare our solution –BRaIn– against eight relevant baselines to place our work in the literature. Through our experiments, we answer three research questions as follows:

- **RQ₁**: (a) How does BRaIn perform in localizing software bugs? (b) Does BRaIn enhance the localization of bugs that require changing multiple documents? (c) Can it improve the localization of bugs that are reported poorly?
- **RQ₂**: How do IRF-based query expansion and document ranking contribute to the performance of BRaIn?
- **RQ₃**: Can BRaIn outperform the relevant baseline techniques in bug localization?

A. Dataset Construction

In our experiment, we used the Bench4BL [21], a comprehensive benchmark dataset that contains 10,017 bug reports from 51 open-source systems, covering a total of 695 software

versions. Our initial assessment revealed that bug reports from the older systems lacked crucial versioning information, making them unsuitable for our study. Additionally, we could not accurately link some bug reports to their corresponding buggy code within the code repositories. Hence, we excluded these bug reports from our dataset. We also found bug reports containing only stack traces without accompanying any textual descriptions of their bugs. We identified those bug reports using regular expressions [10] and excluded them from the dataset. Following these refinement steps, our final dataset comprised 4,683 bug reports from 42 different systems spanning across 684 versions. Table III-a summarizes our curated dataset.

To conduct our experiments and compare with deep learning-based baseline techniques, we used an optimal 80:20 dataset split [65]. This split was done chronologically within each system to ensure that the training set consists of the older 80% of the data, while the test set contains the newest 20% to imitate a real-world scenario. Table III-b provides a summary of our training and test datasets.

B. Evaluation Metrics

Mean Average Precision (MAP): Precision@K indicates the precision for each instance of a buggy source document in the ranked list. Average Precision computes the average Precision@K for all buggy documents in relation to a specific search query. Consequently, Mean Average Precision (MAP) is obtained by averaging the Average Precision values across all queries (Q) within a dataset.

$$AP@K = \frac{1}{|D|} \sum_{k=1}^K P_k \times B_k \quad \Bigg| \quad MAP = \frac{1}{|Q|} \sum_{q=1}^Q AP@K_q$$

Here, $AP@K$ computes average precision for top- K results, where P_k is precision at position k and B_k indicates if item k is buggy (1) or not (0). MAP averages this across all queries q in dataset Q , with D being the ground truth documents.

Mean Reciprocal Rank (MRR): Reciprocal Rank (RR) refers to the rank of the first relevant result retrieved by a technique. It is defined as the reciprocal of the rank of the first relevant source document within the ranked list for each query.

$$RR_q = \frac{1}{\text{Rank of First Relevant Item}} \quad \Bigg| \quad MRR = \frac{1}{|Q|} \sum_{q=1}^{|Q|} RR_q$$

Here, MRR averages the Reciprocal Ranks (RR_q) across all queries q in set Q , where RR_q is the Reciprocal Rank for query q .

HIT@K: HIT@K [35] measures the proportion of queries for which a technique retrieves at least one relevant document among the top- K results. Higher HIT@K values indicate better performance in bug localization techniques.

$$HIT@K = \frac{1}{|Q|} \sum_{q=1}^{|Q|} \begin{cases} 1, & r_q \in \mathcal{G} \\ 0, & \text{otherwise} \end{cases}$$

Here, r_q returns 1 if query q has a ground truth item in the top- K results (0 otherwise), where Q is the set of all queries.

TABLE IV: Performance of BRaIn

Techniques	MAP	MRR	HIT@1	HIT@5	HIT@10
Baseline VSM	0.484	0.513	0.413	0.647	0.732
BRaIn (LLaMA)	0.534	0.568	0.470	0.701	0.766
BRaIn (Mistral)	0.537	0.571	0.469	0.709	0.781
BRaIn (Qwen)	0.492	0.523	0.411	0.678	0.755

C. Selection of LLM

We select three Large Language Models (LLM) to design our techniques and conduct our experiments. We adopt three important criteria to select the models: (a) they should be open-source instruction models, (b) they need to have a quantized version to reduce computational demand, and (c) they should have similar numbers of parameters to allow for fair comparisons. Based on these criteria, we chose LLaMA-3 8B Instruct [54], Mistral v0.3 7B Instruct [34], and Qwen 1.5 7B Chat [58]. They were the top models from instruct category on the Huggingface Open LLM leaderboard [66] during July 2024. We use the 8 bit quantized GPTQ versions of these models [67] to achieve computational efficiency and leverage vLLM [60] for parallelization. We conducted the experiments on Nvidia V100 GPU-enabled machines with 16 GB vRAM in a cluster computing environment.

D. Evaluating BRaIn

Answering RQ_1 - Performance of BRaIn: We evaluate the performance of *BRaIn* using Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), and HIT@K against top 1, 5, and 10 results. Table IV summarizes our performance details.

From Table IV, we see that our proposed technique performs well in detecting the software bugs. BRaIn, powered by Mistral exhibits strong performance, with a Mean Average Precision (MAP) of 0.537. This indicates BRaIn’s ability to rank the relevant documents (a.k.a., buggy source documents) higher than the irrelevant ones. Our technique achieves a Mean Reciprocal Rank (MRR) of 0.570 suggesting that the first relevant document is found within the top two positions. BRaIn (Mistral)’s HIT@1 score of 0.469 shows that, for nearly 47% of bug reports, the most relevant document appears at the top position. BRaIn (Mistral) also performs well in HIT@5 and HIT@10, with approximately 71% and 78% of bug reports having at least one relevant buggy document found within the top 5 and top 10 positions, respectively. BRaIn (LLaMA) delivers nearly comparable results to BRaIn (Mistral), trailing by 1.9% in HIT@10. Although BRaIn (Qwen) demonstrates decent performance, it lags behind both BRaIn (Mistral) and BRaIn (LLaMA) in all metrics. It achieves MAP and MRR scores of 0.492 and 0.523, which are about 9.1% lower than BRaIn (Mistral)’s best performance in each metric.

According to our investigation, some bugs trigger changes to a single document during bug resolution, whereas others trigger changes to multiple documents. We thus evaluate BRaIn’s performance in localizing bugs that warrant changes across multiple source documents. Table V shows the performance of BRaIn, powered by three different LLMs, in terms of

TABLE V: Performance of BRaIn against multi-document bugs

Changed Documents	Bug Report Count	Baseline			LLaMA			Mistral			Qwen		
		MAP	MRR	HIT@10	MAP	MRR	HIT@10	MAP	MRR	HIT@10	MAP	MRR	HIT@10
1	1,949	0.474	0.474	0.690	0.535	0.535	0.726	0.542	0.542	0.739	0.483	0.483	0.712
2	1,436	0.528	0.573	0.767	0.577	0.628	0.801	0.565	0.618	0.820	0.529	0.578	0.792
3	525	0.462	0.518	0.743	0.493	0.554	0.781	0.511	0.573	0.789	0.469	0.528	0.758
4 _≥	773	0.445	0.501	0.765	0.480	0.551	0.793	0.492	0.554	0.811	0.461	0.520	0.794

MAP, MRR, and HIT@10. We grouped bugs from our dataset into four categories based on the number of their changed documents: 1, 2, 3, and 4 or more. Our findings show that BRaIn performs strongest when paired with Mistral. For the 1,949 bugs requiring changes to a single document, BRaIn (Mistral) achieves a MAP of 0.542, representing a significant 14.3% improvement over the baseline counterpart. The improvements extend to other metrics, with a 14.3% increase in MRR and 7.1% in HIT@10. BRaIn (Mistral) also excels in resolving bugs that require multiple document changes, achieving improvements of 7.0-10.6% in MAP, 7.9-10.6% in MRR, and 6.0-6.9% in HIT@10. The other variants of BRaIn also outperform the baseline in localizing bugs that require multiple document changes, achieving improvements of up to 9.2% in MAP, 10.0% in MRR, and 5.1% in HIT@10.

We also investigate how BRaIn performs in localizing bugs where the bug reports could be of low quality (Fig. 3-a). According to existing literature [11], low-quality bug reports lack sufficient information and provide queries that cannot retrieve at least one relevant result within their top 10 positions. In our dataset, we identified 1,101 bug reports that fall into this category. Of these, 581 bug reports (i.e., queries) do not contain any ground truth within their top 50 results returned by Elasticsearch. These bug reports were not considered, which leaves us with 520 low-quality bug reports for our analysis. Our findings demonstrate BRaIn’s promising results even with the low-quality reports. BRaIn (Mistral) emerges as the top performer, successfully localizing 268 bug reports (51.5% of low-quality bug reports) within the top 10 results. BRaIn (LLaMA) and BRaIn (Qwen) follow, identifying 225 and 198 reports, respectively. Notably, all three models identified 130 bugs, with BRaIn (Mistral) uniquely localizing an additional 70 bugs, followed by BRaIn (LLaMA) and BRaIn (Mistral) with 30 and 23 bugs. We also assess BRaIn’s capability to detect the first relevant document (a.k.a., buggy source document), where 130 bug reports from above were considered (Fig. 3-b). Interestingly, BRaIn (Qwen) outperformed the other variants in this metric, localizing 26.9% of the bugs at the top positions, followed by BRaIn (LLaMA) at 18.5% and BRaIn (Mistral) at 10%. All the findings above suggest BRaIn’s ability to analyze, enhance, and localize bug reports, even with low-quality reports.

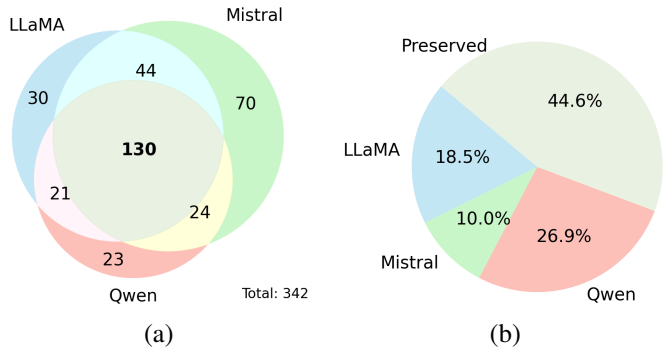


Fig. 3: Performance of BRaIn with Low Quality Bug Reports

RQ1 Summary: BRaIn significantly improves bug localization, particularly with Mistral, reaching a high MAP score of 0.537. This performance is due to BRaIn’s effective handling of bug reports with up to $\approx 11\%$ multiple changed documents. Moreover, BRaIn demonstrates an impressive performance with limited information, successfully localizing $\approx 52\%$ of low-quality bug reports within the top 10 results where the baseline failed.

TABLE VI: Impact of Query Expansion and Scoring

BRaIn Components	BRaIn (LLaMA)		BRaIn (Mistral)		BRaIn (Qwen)	
	MAP	MRR	MAP	MRR	MAP	MRR
Expansion + Reranking	0.534	0.568	0.537	0.571	0.492	0.523
Expansion + No Reranking	0.498	0.568	0.498	0.569	0.496	0.567
No Expansion + Reranking	0.518	0.574	0.520	0.573	0.489	0.529

Answering RQ₂ - Contribution of Query Expansion and Document reranking: BRaIn leverages Intelligent Relevant Feedback (IRF) to expand its original query and rerank the documents. Query expansion can improve bug localization by adding relevant keywords to an original query. Similarly, incorporating contextual understanding into document scoring can help go beyond just textual and semantic matching during document ranking. We examine the contribution of these two components (Table VI) to BRaIn’s performance as follows.

To determine the impact of query expansion in isolation, we evaluated BRaIn’s performance without the reranking component (Table VI). Interestingly, all BRaIn variants achieved similar MAP scores of around 0.49, falling short of opti-

TABLE VII: Comparison Between BRaIn and Baseline Techniques

(a) Comparison with Non-ML Baselines

Metrics	Traditional IR			Relevance Feedback		
	Baseline	BLUiR	Blizzard	Rocchio	Sysman-SCP	BRaIn
MAP	0.484	0.450	0.506	0.489	0.472	0.537
MRR	0.513	0.471	0.536	0.558	0.541	0.571
HIT@10	0.732	0.696	0.758	0.765	0.753	0.781

(b) Comparison with ML Baselines

Metrics	Machine Learning			
	DNNLOC	RLocator	NextBug	BRaIn
MAP	0.283	0.488	0.469	0.531
MRR	0.296	0.561	0.540	0.564
HIT@10	0.518	0.735	0.743	0.771

mal performance. For example, BRaIn (LLaMA) and BRaIn (Mistral) saw MAP scores decrease by 7.2% and 7.8%, respectively, while their MRR scores remained relatively stable. BRaIn (Qwen), on the other hand, showed a 7.8% increase in MRR alongside a slight improvement in MAP. However, all BRaIn variants outperformed the baseline by 10.5-10.9% in MRR and 2.4-2.9% in MAP.

In contrast, when reranking was applied isolately, MAP scores for BRaIn (LLaMA), BRaIn (Mistral), and BRaIn (Qwen) dropped by 3.0%, 3.2%, and 6.1%, respectively, while their MRR scores remained consistent. Despite these decreases, all variants showed improvements of 3.1-11% in MRR and 1.0-7.1% in MAP over the baseline.

These findings show that Intelligent Relevance Feedback improves the individual components of query expansion and reranking. However, they also underscore the importance of the synergy between these components in achieving optimal performance.

RQ2 Summary: Query expansion and reranking individually decrease MAP and MRR by 3.0-7.8%, yet both improve bug localization performance over the baseline by 2.4-11% in these metrics. Their individual results highlight the contribution of Intelligent Relevance Feedback and underscore the importance of a synergistic combination in BRaIn.

Answering RQ₃ - Comparison with Baseline Techniques:

To place our work in the literature, we compare BRaIn with relevant baseline techniques in terms of their MAP, MRR, and HIT@10. Given our methodology, we choose two types of baseline techniques – IR methods [10], [29], [35] and deep learning based methods [36]–[38]. For comparison with baselines, we use BRaIn (Mistral) in our experiments since it is the best-performing variant of BRaIn.

To replicate the traditional IR-based Baseline VSM, we index all source documents of a repository and use bug reports (title + description) as queries. These queries are executed with the Elasticsearch [40], which retrieves relevant documents using the BM25 algorithm [39] and Boolean search, with default parameters for k and b . Other traditional approaches from literature– BLUiR [35] and Blizzard [10]– use structured information from bug reports and source code for bug localization. BLUiR calculates suspiciousness scores using class names, method names, variable names, comments, and bug report elements (title, description), combining multiple

searches into an overall score. Blizzard categorizes bug reports into three types and constructs text graphs from these reports to generate queries and retrieve relevant buggy source documents. For both approaches, we employ Apache Lucene [68] for retrieval. We replicated these methods by adapting them from Bench4BL repository [21] and Blizzard’s replication package [69] from the authors. We compare BRaIn with these IR-based techniques to assess the performance of our technique against established models in the field.

Relevance feedback-based techniques like Rocchio [70] and the Spatial Code Proximity (SCP) model [29] aim to enhance bug localization by refining queries based on the results of initial searches. Rocchio is a widely-used relevance feedback technique for information retrieval [70]. We leverage relevance feedback to reformulate queries and Apache Lucene to execute the queries and retrieve the documents. Our reformulated queries were optimized using α , β , and γ parameters [70]. Similarly, we implemented Sisman et al.’s SCP model [29] to reformulate queries based on term proximity within source code. It prioritizes terms that frequently co-occur within the same method or class, using the best parameters w , x , and y suggested by the authors. We compare these techniques to our approach to highlight the importance of contextual understanding during relevance feedback of search queries.

Since Machine Learning (ML) techniques can capture complex patterns in data using non-linear relationships, we compare BRaIn against three ML-based techniques– DNNLOC [36], NextBug [37], and RLocator [38]. DNNLOC combines multiple features– rVSM score [71] for bug report-source code similarity, class name similarity, collaborative filtering, and bug report recency and frequency—and uses a neural network to predict suspiciousness scores to rank documents. NextBug employs Word2Vec [72] embeddings to capture semantic relations between bug reports and source code and thus to localize the buggy documents. In our experiments, we substituted Word2Vec with CodeT5 embeddings [73] to capture more nuanced text-level semantic associations, as opposed to token-level. RLocator is a recent deep-learning technique that employs a reinforcement learning model, framed as a Markov Decision Process, to optimize ranking of buggy documents. We replicated DNNLOC and NextBug by following the respective authors’ approaches and replicate RLocator using the authors’ provided replication package on Zenodo [74].

Table VII-a summarizes our comparison details with the baseline techniques. Among traditional IR-based approaches, Blizzard achieves a MAP score of 0.506, while Baseline VSM and BLUiR scores are 0.484 and 0.450. BRaIn out-

performs them with a MAP of 0.537, achieving a maximum improvement of 19.3% over these techniques. Similarly, BRaIn achieves notable gains in MRR and HIT@10, with increases of up to 17.5% and 12.2%. Among the IR based approaches that leverage relevance feedback, Rocchio’s algorithm achieves a MAP of 0.489, slightly above the baseline, while Sysman-SCP falls short by 2.5%. BRaIn again leads here with the improvements in MAP, MRR, and HIT@10 of 13.8%, 5.5%, and 3.7%, respectively. These results underscore the advantages of Intelligent Relevance Feedback, which uses contextual understanding over traditional techniques based on textual relevance for bug localization.

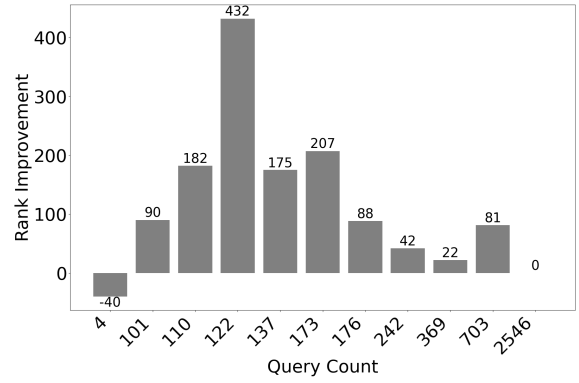
As shown in VII-b, BRaIn also outperforms the machine learning techniques that require training. We evaluated both BRaIn and the baseline techniques on the test set only to ensure a fair comparison. It should be noted that old bug reports and their corresponding code were used for training and the recent bugs and their corresponding code were used for testing. DNNLOC performs significantly lower with a MAP of 0.283, 87.6% lower than BRaIn’s optimal score of 0.531. In comparison, RLocator and NextBug achieve MAP scores of 0.488 and 0.469, with BRaIn outperforming them by 8.8% and 13.2%, respectively. Similar improvements are observed for other metrics, with BRaIn showing 4.4–89.5% improvements in MRR and 3.7–48.8% in HIT@10. Such performance underscores the superiority of BRaIn’s performance with intelligent relevance feedback (IRF) compared to baseline techniques.

TABLE VIII: Statistical Test: BRaIn vs. Blizzard

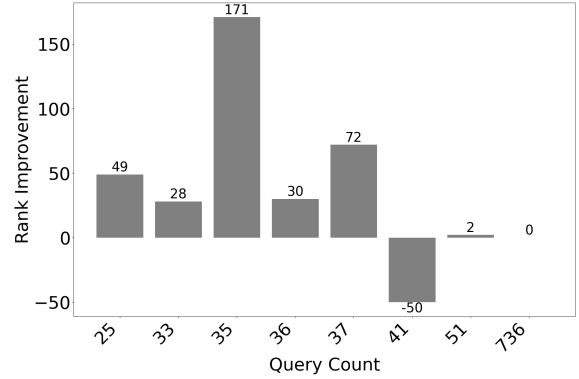
Evaluation Point	p-value	Effect Size (Cliff’s δ)
Top-1	0.0023 **	Medium (0.41)
Top-5	0.0015 **	Large (0.66)
Top-10	0.0008 ***	Large (0.82)

*=statistical significance

Finding the first buggy document is very important during bug localization [29]. We further investigate how BRaIn performs in such a case. We chose Blizzard for this investigation as it is the best-performing model against BRaIn in our experiments. Fig. 4-a compares BRaIn and Blizzard by analyzing the difference in ranks for each query within the top 10 results. A positive value indicates that BRaIn found the first ground truth at a better rank than Blizzard, while a negative value suggests the opposite. For 122 bug reports, the large difference of 432 indicates that BRaIn identified the first buggy documents more often than Blizzard. In contrast, Blizzard only outperformed BRaIn for 4 bug reports. For the set of 2,546 bug reports, there are two distinct possibilities: Either (a) the rank difference between the two techniques was 0, or (b) neither technique identified a buggy document within the top-10 results. We extend our analysis to the 1,101 low-quality bug reports discussed in RQ_1 . Here we also see BRaIn’s dominance in rank improvement over Blizzard for 80.34% low-quality bug reports (Fig. 4-b). These results strongly indicate the superiority of our approach. To further validate our findings against Blizzard, we conducted non-parametric



(a) Full Dataset (4,683 Bug Reports)



(b) Low Quality Bug Reports (1,101 Bug Reports)

Fig. 4: Rank Improvement: BRaIn vs Blizzard

statistical tests –Mann-Whitney Wilcoxon and Cliff’s δ [75]– and compared BRaIn to Blizzard in identifying the first buggy document on the entire dataset. The tests in Table VIII show p -values < 0.05 for the top-1, 5, and 10 results, with medium to large effect sizes (i.e., $0.41 \leq \delta \leq 0.82$). Thus, it confirms BRaIn’s consistent ability to detect the first buggy document more effectively than its closest competitor.

RQ3 Summary: BRaIn outperforms traditional, relevance feedback-based, and ML-based baseline techniques, achieving an 87.6% improvement in MAP. This demonstrates BRaIn’s ability to localize relevant buggy documents at the top positions using Intelligent Relevance Feedback (IRF), surpassing relevance feedback-based techniques by 3.7–13.8% across various metrics. Statistical significance tests further confirm BRaIn’s superiority.

V. RELATED WORK

A. IR based Bug Localization

Bug localization techniques can broadly be classified into two main groups: spectra-based and information retrieval (IR)-based approaches [14]. Spectra-based methods use program execution traces and test methods to localize bugs, making them complex and expensive [12], [13]. In contrast, IR-based

techniques rely on textual overlap between bug reports and source code to localize the bugs.

Traditional IR-based bug localization methods that leverage the vector space model (VSM) [76], have been enhanced by integrating additional contexts, such as bug report history, code modifications, and version history [23], [77], [78]. For instance, Saha et al. [35] leverage bug report and source code structures, capture eight components from the code and bug reports, and perform eight pairwise searches using a sophisticated retrieval technique, Indri [79]. On the other hand, BugLocator [71] combines a modified VSM (rVSM [71]) score with previous bug fix history to improve bug localization. AmaLgam [19] integrates BLUiR, BugLocator, and version history to better detect buggy documents. AmaLgam+ [80] further incorporates stack traces and bug reporter history, refining bug localization across five ranking components. While advanced, computationally expensive methodologies such as LSI or LDA [71], [81] are available, their bug localization effectiveness is similar to that of more basic methods [21].

In our work, we use a VSM-based retrieval that is enhanced by Boolean search (e.g., Elasticsearch). However, it was complemented by contextual understanding of bugs and Intelligent Relevance Feedback, leveraging the capabilities of LLMs.

B. Query Reformulation

Poorly constructed queries from software bug reports can significantly hinder IR-based bug localization [11], [82]. To tackle this issue, researchers have developed query reformulation techniques that can improve search queries by incorporating better terms or eliminating unnecessary ones. For instance, Refoqus [29] uses query characteristics and machine learning to recommend strategies like query reduction or expansion for a given query. Graph-based methods analyze semantic and syntactic relationships within bug reports to identify key terms. Rahman and Roy [10] create text graphs to collect important terms based on three different bug types to improve query reformulation. The authors later demonstrate the presence of optimal queries in bug reports by employing Genetic algorithms [11]. Relevance feedback-based approaches iteratively refine queries based on their search results. Gay et al. [29] used Rocchio’s algorithm to improve queries with developer feedback, while Sisman et al. [30] expanded queries by selecting terms from the top-ranked documents using Spatial Code Proximity (SCP), without using any explicit relevance feedback.

These approaches rely on statistical properties or co-occurrence relations to reformulate an original query without meaningful knowledge of source code or bug reports. Our approach digs deeper to select relevant source code by contextually understanding bug reports to formulate queries for better bug resolution.

C. Deep Learning and Bug Localization

Recent advancements in deep learning have encouraged its applications in bug localization. DNNLOC [36], a seminal work on this topic, identifies buggy documents by learning

from multiple text-based features and metadata (e.g., rVSM score [71], class name similarity, bug report recency). However, its reliance on features like bug fixing recency can limit its application [71]. A recent technique, FBL-BERT [83] uses a BERT-based model, ColBERT [84], for document scoring with late interaction. However, it relies on changesets for resolution, which are difficult to track in large, fast-changing projects. Another recent technique, RLocator [38], aims to optimize ranking metrics in the bug localization process. It formulates bug localization as a Markov Decision Process (MDP) and employs reinforcement learning (RL) to guide its decisions. Other approaches like TRANP-CNN [85] and CooBa [86] use Convolutional Neural Networks (CNN) and Graph Convolutional Networks (GCN) to improve cross-project bug localization. However, these methods may suffer from a lack of scalability when handling large volumes of documents.

In contrast, our technique deals with a limited set of documents for scalability, retrieved by a scalable tool – Elasticsearch [40]. By combining efficient IR-based filtering with the contextual depth of language models, BRAIn accurately narrows down and ranks buggy documents, improving bug localization.

VI. THREATS TO VALIDITY

Threats to internal validity concern experimental errors and biases. Replication of existing baselines poses such a threat. We mitigated this by using replication packages from the original authors (Blizzard [69], RLocator [74]) and from Bench4BL [21] (BLUiR [35]). Due to Indri’s obsolescence, we substituted it with Lucene in the BLUiR replication. On the other hand, we replicated DNNLOC [36], NextBug [37], and Sysman-SCP [29] adhering strictly to the original authors’ settings and parameters. To minimize bias, we tested on two distinct datasets and found only a negligible difference compared to baseline performances.

Threats to external validity relate to generalizability. While BRAIn was evaluated only on Java code, the underlying models (e.g., LLaMa [54]) are designed to adapt to various programming languages, potentially mitigating this limitation.

Threats to construct validity concern the appropriateness of our evaluation metrics. We employed widely used metrics such as Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), and HIT@K, which were commonly used in existing literature on bug localization [10], [35], [87] and Information Retrieval studies [88], [89]. Therefore, this choice of metrics minimizes threats to construct validity.

Finally, we used 20 bug reports from our dataset to optimize prompts with LLaMA. Since they are a part of our experimental dataset, it could introduce bias. We repeated a limited experiment and found similar performance to the reported ones in the paper. Thus, any relevant threat of bias is minimal to none.

VII. CONCLUSION AND FUTURE WORK

Software bugs consume resources and cause financial losses [4]. Resolving these challenges has been a major focus, with

researchers working on solutions for decades. In this paper, we introduced BRAIn to address the contextual gaps between bug reports and source codes by assessing the relevance between bug reports and code with Large Language Models (LLM). Our proposed technique leverages the LLM’s feedback (a.k.a., Intelligent Relevance Feedback) to improve bug localization by reformulating queries and reranking source documents. We evaluated BRAIn’s performance using three widely used metrics in bug localization: Mean Average Precision (MAP), Mean Reciprocal Rank (MRR), and HIT@K. BRAIn achieved substantial improvements, showing increases of 87.6% in MAP, 89.5% in MRR, and 48.8% in HIT@K compared to baseline techniques. Furthermore, BRAIn detected $\approx 52\%$ more low-quality bugs within the top 10 ranks relative to the baseline VSM.

Building on these results, in the future we aim to localize bugs using in few-shot settings [52] and at a much finer level (i.e., method, line) to assist developers in pinpointing the exact location of bug inducing code.

REFERENCES

- [1] C. for IT Software Quality, “Cpsq 2020 report,” IT-CISQ, Tech. Rep., 2020. [Online]. Available: <https://www.it-cisq.org/cisq-files/pdf/CPSQ-2020-report.pdf>
- [2] F. T. Council. (2023) Costly code: The price of software errors. [Online]. Available: <https://shorturl.at/2j7wu>
- [3] D. Weston, “Helping our customers through the crowdstrike outage - the official microsoft blog,” Jul. 2024. [Online]. Available: <https://shorturl.at/3OCxV>
- [4] L. K. Wee, “Here comes the wave of insurance claims for the crowdstrike outage,” Jul. 2024. [Online]. Available: <https://shorturl.at/5Y1jQ>
- [5] W. Zou, D. Lo, Z. Chen, X. Xia, Y. Feng, and B. Xu, “How practitioners perceive automated bug report management techniques,” *IEEE Transactions on Software Engineering*, vol. 46, no. 8, pp. 836–862, 2018.
- [6] D. H. O’Dell, “The debugging mindset: Understanding the psychology of learning strategies leads to effective problem-solving skills.” *Queue*, vol. 15, no. 1, pp. 71–90, 2017.
- [7] DevOps. (2024) Survey: Fixing bugs stealing time from development. [Online]. Available: <https://shorturl.at/Fj8sB>
- [8] T. Britton, L. Jeng, G. Carver, P. Cheak, and T. Katzenellenbogen, “Reversible debugging software,” *Judge Bus. School, Univ. Cambridge, Cambridge, UK, Tech. Rep.*, vol. 229, 2013.
- [9] G. W. Furnas, T. K. Landauer, L. M. Gomez, and S. T. Dumais, “The vocabulary problem in human-system communication,” *Communications of the ACM*, vol. 30, no. 11, pp. 964–971, 1987.
- [10] M. M. Rahman and C. K. Roy, “Improving ir-based bug localization with context-aware query reformulation,” in *Proceedings of the 2018 26th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*, 2018, pp. 621–632.
- [11] M. M. Rahman, F. Khomh, S. Yeasmin, and C. K. Roy, “The forgotten role of search queries in ir-based bug localization: an empirical study,” *Empirical Software Engineering*, vol. 26, no. 6, p. 116, 2021.
- [12] L. Moreno, J. J. Treadway, A. Marcus, and W. Shen, “On the use of stack traces to improve text retrieval-based bug localization,” in *2014 IEEE International Conference on Software Maintenance and Evolution*, 2014, pp. 151–160.
- [13] Q. Wang, C. Parnin, and A. Orso, “Evaluating the usefulness of ir-based fault localization techniques,” in *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, ser. ISSTA 2015. New York, NY, USA: Association for Computing Machinery, 2015, p. 1–11. [Online]. Available: <https://doi.org/10.1145/2771783.2771797>
- [14] —, “Evaluating the usefulness of ir-based fault localization techniques,” in *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, 2015, pp. 1–11.
- [15] A. T. Nguyen, T. T. Nguyen, J. Al-Kofahi, H. V. Nguyen, and T. N. Nguyen, “A topic-based approach for narrowing the search space of buggy files from a bug report,” in *2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011)*. IEEE, 2011, pp. 263–272.
- [16] T.-D. B. Le, R. J. Oentaryo, and D. Lo, “Information retrieval and spectrum based bug localization: Better together,” in *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*, 2015, pp. 579–590.
- [17] A. N. Lam, A. T. Nguyen, H. A. Nguyen, and T. N. Nguyen, “Bug localization with combination of deep learning and information retrieval,” in *2017 IEEE/ACM 25th International Conference on Program Comprehension (ICPC)*. IEEE, 2017, pp. 218–229.
- [18] X. Wang, C. Macdonald, and I. Ounis, “Deep reinforced query reformulation for information retrieval.” *arXiv preprint arXiv:2007.07987*, 2020.
- [19] S. Wang and D. Lo, “Version history, similar report, and structure: Putting them together for improved bug localization,” in *Proceedings of the 22nd International Conference on Program Comprehension*, 2014, pp. 53–63.
- [20] K. C. Youm, J. Ahn, J. Kim, and E. Lee, “Bug localization based on code change histories and bug reports,” in *2015 Asia-Pacific Software Engineering Conference (APSEC)*, 2015, pp. 190–197.
- [21] J. Lee, D. Kim, T. F. Bissyandé, W. Jung, and Y. Le Traon, “Bench4bl: reproducibility study on the performance of ir-based bug localization,” in *Proceedings of the 27th ACM SIGSOFT international symposium on software testing and analysis*, 2018, pp. 61–72.
- [22] O. Chaparro, J. M. Florez, and A. Marcus, “Using observed behavior to reformulate queries during text retrieval-based bug localization,” in *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 2017, pp. 376–387.
- [23] B. Sisman and A. C. Kak, “Incorporating version histories in information retrieval based bug localization,” in *2012 9th IEEE working conference on mining software repositories (MSR)*. IEEE, 2012, pp. 50–59.
- [24] C. Li, Y. Sun, B. He, L. Wang, K. Hui, A. Yates, L. Sun, and J. Xu, “Nprf: A neural pseudo relevance feedback framework for ad-hoc information retrieval,” Jan. 2018. [Online]. Available: <https://arxiv.org/abs/1810.12936>
- [25] S. Yu, D. Cai, J.-R. Wen, and W.-Y. Ma, “Improving pseudo-relevance feedback in web information retrieval using web page segmentation,” in *Proceedings of the 12th international conference on World Wide Web*, 2003, pp. 11–18.
- [26] J. Wang, M. Pan, T. He, X. Huang, X. Wang, and X. Tu, “A pseudo-relevance feedback framework combining relevance matching and semantic matching for information retrieval,” *Information Processing & Management*, vol. 57, no. 6, p. 102342, 2020.
- [27] M. Kim and E. Lee, “A novel approach to automatic query reformulation for ir-based bug localization,” in *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing*, ser. SAC ’19. New York, NY, USA: Association for Computing Machinery, 2019, p. 1752–1759. [Online]. Available: <https://doi.org/10.1145/3297280.3297451>
- [28] G. Gay, S. Haiduc, A. Marcus, and T. Menzies, “On the use of relevance feedback in ir-based concept location,” in *2009 IEEE International Conference on Software Maintenance*, 2009, pp. 351–360.
- [29] S. Haiduc, G. Bavota, A. Marcus, R. Oliveto, A. De Lucia, and T. Menzies, “Automatic query reformulations for text retrieval in software engineering,” in *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 2013, pp. 842–851.
- [30] B. Sisman and A. C. Kak, “Assisting code search with automatic query reformulation for bug localization,” in *2013 10th Working Conference on Mining Software Repositories (MSR)*, 2013, pp. 309–318.
- [31] M. Allamanis, E. T. Barr, P. Devanbu, and C. Sutton, “A survey of machine learning for big code and naturalness,” *ACM Computing Surveys (CSUR)*, vol. 51, no. 4, pp. 1–37, 2018.
- [32] F. Deissenboeck and M. Pizka, “Concise and consistent naming,” *Software Quality Journal*, vol. 14, pp. 261–282, 2006.
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [34] M. Team, “Mistral: A new approach to language models,” 2023. [Online]. Available: <https://mistral.ai>
- [35] R. K. Saha, M. Lease, S. Khurshid, and D. E. Perry, “Improving bug localization using structured information retrieval,” in *2013 28th*

- IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2013, pp. 345–355.
- [36] A. N. Lam, A. T. Nguyen, H. A. Nguyen, and T. N. Nguyen, “Bug localization with combination of deep learning and information retrieval,” in *2017 IEEE/ACM 25th International Conference on Program Comprehension (ICPC)*, 2017, pp. 218–229.
- [37] “Combining word embedding with information retrieval to recommend similar bug reports.” IEEE Computer Society, 12 2016, pp. 127–137.
- [38] P. Chakraborty, M. Alfadel, and M. Nagappan, “Rlocator: Reinforcement learning for bug localization,” *IEEE Transactions on Software Engineering*, 2024.
- [39] S. E. Robertson, S. Walker, S. Jones, M. M. Hancock-Beaulieu, M. Gatford *et al.*, “Okapi at trec-3,” *Nist Special Publication Sp*, vol. 109, p. 109, 1995.
- [40] “Elasticsearch,” Elastic. [Online]. Available: <https://www.elastic.co/elasticsearch/>
- [41] Q. Gu, “Llm-based code generation method for golang compiler testing,” in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2023, pp. 2201–2203.
- [42] L. Lian, B. Li, A. Yala, and T. Darrell, “Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models,” *arXiv preprint arXiv:2305.13655*, 2023.
- [43] Y. Fathullah, C. Wu, E. Lakomkin, J. Jia, Y. Shangguan, K. Li, J. Guo, W. Xiong, J. Mahadeokar, O. Kalinli *et al.*, “Prompting large language models with speech recognition abilities,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 13 351–13 355.
- [44] A. Webson and E. Pavlick, “Do prompt-based models really understand the meaning of their prompts?” *arXiv preprint arXiv:2109.01247*, 2021.
- [45] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, “Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing,” *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.
- [46] M. Mukhtadir, Golam, “A brief history of prompt: Leveraging language models. (through advanced prompting),” Sep. 2023. [Online]. Available: <https://arxiv.org/abs/2310.04438>
- [47] B. Meskó, “Prompt engineering as an important emerging skill for medical professionals: tutorial,” *Journal of medical Internet research*, vol. 25, p. e50638, 2023.
- [48] B. Chen, Z. Zhang, N. Langrené, and S. Zhu, “Unleashing the potential of prompt engineering in large language models: a comprehensive review,” *arXiv preprint arXiv:2310.14735*, 2023.
- [49] S. Brade, B. Wang, M. Sousa, S. Oore, and T. Grossman, “Promptify: Text-to-image generation through interactive prompt exploration with large language models,” in *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 2023, pp. 1–14.
- [50] S. Feng and C. Chen, “Prompting is all you need: Automated android bug replay with large language models,” in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–13.
- [51] J. Shin, C. Tang, T. Mohati, M. Nayebi, S. Wang, and H. Hemmati, “Prompt engineering or fine tuning: An empirical assessment of large language models in automated software engineering tasks,” *arXiv preprint arXiv:2310.10508*, 2023.
- [52] S. Schulhoff, M. Ilie, N. Balepur, K. Kahadze, A. Liu, C. Si, Y. Li, A. Gupta, H. Han, S. Schulhoff *et al.*, “The prompt report: A systematic survey of prompting techniques,” *arXiv preprint arXiv:2406.06608*, 2024.
- [53] T. Schnabel and J. Neville, “Prompts as programs: A structure-aware approach to efficient compile-time prompt optimization,” *arXiv preprint arXiv:2404.02319*, 2024.
- [54] H. Touvron, A. Bosselut, K. Sinha, and *et al.*, “Llama: Open and efficient foundation language models,” 2023. [Online]. Available: <https://arxiv.org/abs/2302.13971>
- [55] H. Luo, L. Jiang, Y. Belinkov, and J. Glass, “Improving neural language models by segmenting, attending, and predicting the future,” *arXiv preprint arXiv:1906.01702*, 2019.
- [56] V. Srinivasan and T. Reps, “An improved algorithm for slicing machine code,” *ACM SIGPLAN Notices*, vol. 51, no. 10, pp. 378–393, 2016.
- [57] C. Fischer, “Javaparser,” 2019. [Online]. Available: <https://github.com/javaparser/javaparser>
- [58] Q. Team, “Qwen: A high-performance language model,” 2023. [Online]. Available: <https://huggingface.co/Qwen>
- [59] T. Wolf, L. Debut, V. Sanh, J. Chaumond, and Delangue, “Huggingface’s transformers: State-of-the-art natural language processing,” *arXiv preprint arXiv:1910.03771*, 2019. [Online]. Available: <https://huggingface.co/docs/transformers/index>
- [60] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and I. Stoica, “Efficient memory management for large language model serving with pagedattention,” in *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [61] M. M. Rahman and C. K. Roy, “Improved query reformulation for concept location using coderank and document structures,” in *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2017, pp. 428–439.
- [62] “Javalang: Pure python java parser and tools.” [Online]. Available: <https://github.com/c2nes/javalang>
- [63] S. Brin and L. Page, “The anatomy of a large-scale hypertextual web search engine,” *Computer networks and ISDN systems*, vol. 30, no. 1-7, pp. 107–117, 1998.
- [64] M. Franke and J. Degen, “The softmax function: Properties, motivation, and interpretation,” 2023.
- [65] A. Gholamy, V. Kreinovich, and O. Kosheleva, “Why 70/30 or 80/20 relation between training and testing sets: A pedagogical explanation,” *Int. J. Intell. Technol. Appl. Stat.*, vol. 11, no. 2, pp. 105–111, 2018.
- [66] F. Author and S. Author, “Open-llm-leaderboard: From multi-choice to open-style questions for large language models,” *arXiv preprint arXiv:2406.07545*, 2024.
- [67] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, “Gptq: Accurate post-training quantization for generative pre-trained transformers,” *arXiv preprint arXiv:2210.17323*, 2022.
- [68] T. A. S. Foundation, “Apache lucene,” 2021. [Online]. Available: <https://lucene.apache.org/>
- [69] M. Rahman, “Blizzard.” [Online]. Available: <https://github.com/masud-technope/BLIZZARD>
- [70] R. Cummins, N. Language, and I. P. N. Group, “Lecture 7: Relevance feedback and query expansion,” p. 271, 2017. [Online]. Available: <https://www.cl.cam.ac.uk/teaching/1617/InfoRtrv/lecture7-relevance-feedback.pdf?formCode=MG0AV3>
- [71] D. Kim, Y. Tao, S. Kim, and A. Zeller, “Where should we fix this bug? a two-phase recommendation model,” *IEEE transactions on software Engineering*, vol. 39, no. 11, pp. 1597–1610, 2013.
- [72] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” 2013. [Online]. Available: <https://arxiv.org/abs/1310.4546>
- [73] Y. Wang, W. Wang, S. Joty, and S. C. Hoi, “Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation,” *arXiv preprint arXiv:2109.00859*, 2021.
- [74] Zenodo, Jul. 2024. [Online]. Available: <https://doi.org/10.5281/zenodo.11265302>
- [75] T. Barnes, S. C. Moore, and K. Osatuke, *Testing Significance Tests: A Simulation with Cliff’s Delta, t-tests, and Mann-Whitney U*. National Center for Organizational Development, Department of Veteran Affairs, 2018.
- [76] D. Lee, H. Chuang, and K. Seamons, “Document ranking and the vector-space model,” *IEEE Software*, vol. 14, no. 2, pp. 67–75, 1997.
- [77] R. K. Saha, J. Lawall, S. Khurshid, and D. E. Perry, “On the effectiveness of information retrieval based bug localization for c programs,” in *2014 IEEE international conference on software maintenance and evolution*. IEEE, 2014, pp. 161–170.
- [78] M. Wen, R. Wu, and S.-C. Cheung, “Locus: Locating bugs from software changes,” in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, 2016, pp. 262–273.
- [79] T. Strohmaier, D. Metzler, H. R. Turtle, and W. B. Croft, “Indri : A language-model based search engine for complex queries (extended version),” 2005. [Online]. Available: <https://api.semanticscholar.org/CorpusID:18471028>
- [80] S. Wang and D. Lo, “Amalgam+: Composing rich information sources for accurate bug localization,” *Journal of Software: Evolution and Process*, vol. 28, no. 10, pp. 921–942, 2016.
- [81] D. Poshvanyk, Y.-G. Guneuc, A. Marcus, G. Antoniol, and V. Rajlich, “Feature location using probabilistic ranking of methods based on execution scenarios and information retrieval,” *IEEE Transactions on Software Engineering*, vol. 33, no. 6, pp. 420–432, 2007.
- [82] C. Mills, E. Parra, J. Pantuchina, G. Bavota, and S. Haiduc, “On the relationship between bug reports and queries for text retrieval-based bug

- localization,” *Empirical Software Engineering*, vol. 25, pp. 3086–3127, 2020.
- [83] A. Ciborowska and K. Damevski, “Fast changeset-based bug localization with bert,” in *Proceedings of the 44th International Conference on Software Engineering*, 2022, pp. 946–957.
- [84] O. Khattab and M. Zaharia, “Colbert: Efficient and effective passage search via contextualized late interaction over bert,” in *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, 2020, pp. 39–48.
- [85] X. Huo and M. Li, “Enhancing the unified features to locate buggy files by exploiting the sequential nature of source code.” in *IJCAI*, 2017, pp. 1909–1915.
- [86] Z. Zhu, Y. Li, H. Tong, and Y. Wang, “Cooba: Cross-project bug localization via adversarial transfer learning,” in *IJCAI*, 2020.
- [87] M. M. Rahman, C. K. Roy, and D. Lo, “Rack: Automatic api recommendation using crowdsourced knowledge.” *Institute of Electrical and Electronics Engineers (IEEE)*, 5 2016, pp. 349–359.
- [88] D. Harman, *Information retrieval evaluation*. Morgan & Claypool Publishers, 2011.
- [89] L. Carnevali, “Evaluation Measures in Information Retrieval,” <https://www.pinecone.io/learn/offline-evaluation/>.