

# Exploiting spatial code proximity and order for improved source code retrieval for bug localization

Bunyamin Sisman, Shayan A. Akbar and Avinash C. Kak<sup>\*†</sup>

Department of Electrical and Computer Engineering, Purdue University, West Lafayette, IN 47907, USA

## ABSTRACT

Practically all information retrieval based approaches developed to date for automatic bug localization are based on the bag-of-words assumption that ignores any positional and ordering relationships between the terms in a query. In this paper, we argue that bug reports are ill-served by this assumption because such reports frequently contain various types of structural information whose terms must obey certain positional and ordering constraints. It therefore stands to reason that the quality of retrieval for bug localization would improve if these constraints could be taken into account when searching for the most relevant files. In this paper, we demonstrate that such is indeed the case. We show how the well-known Markov Random Field based retrieval framework can be used for taking into account the term-term proximity and ordering relationships in a query vis-à-vis the same relationships in the files of a source-code library to greatly improve the quality of retrieval of the most relevant source files. We have carried out our experimental evaluations on popular large software projects using over 4000 bug reports. The results we present demonstrate unequivocally that the new proposed approach is far superior to the widely used bag-of-words based approaches. Copyright © 2016 John Wiley & Sons, Ltd.

Received 18 April 2016; Revised 12 July 2016; Accepted 18 July 2016

KEY WORDS: bug localization; source code search; term proximity; Markov Random Fields

## 1. INTRODUCTION

Code search plays an important role in software development and maintenance. The tools that are deployed today for code search range all the way from simple command-line functions like ‘grep’ to complex search facilities tailored for the specific needs of the developers. These different types of search facilities are used to locate various parts of a software library for concept location, change impact analysis, traceability link analysis, and so on [1–6]. Our particular interest lies in a class of code search tools that are based on information retrieval (IR) techniques, and our end goal is bug localization. For automatic bug localization, we treat the bug reports as text queries and the corresponding software constructs that should be modified to implement a fix for the bugs as the relevant artifacts that should be retrieved by the search tool [7–9].

Obviously, the success of an IR framework that leverages bug reports for automatic bug localization depends much on how the bug reports are represented vis-à-vis the source code files and other documents in a library. In the widely used *bag-of-words* (BOW) representations for both the queries and the source code files, all positional and ordering relationships between the terms are lost [8–10]. To us, this is tantamount to a serious loss of information considering that a bug report, in general, is a composition of structured and unstructured textual data that frequently includes the following: (i) patches; (ii) stack traces when the software fault throws an exception; (iii) snippets of code; and

<sup>\*</sup>Correspondence to: Avinash C. Kak, Department of Electrical and Computer Engineering, Purdue University, West Lafayette, IN 47907, USA.

<sup>†</sup>E-mail: kak@purdue.edu

(iv) natural language sentences [7, 11]. Patches and stack traces, especially, contain vital inter-term proximity and ordering relationships that ought to be exploited for the purpose of retrieval. For instance, if two terms are proximal to each other in a stack trace, the user would want a source code file containing similar code to have the same two terms in a similar proximal relationship. The work we report in this paper demonstrates that the quality of retrieval improves greatly when a retrieval framework allows for ordering and positional (through proximity) relationships to be taken into account in the retrieval process.

As to how to incorporate ordering and positional relationships in a retrieval framework, we have several possibilities at our disposal that have been examined in the past mostly in the context of retrieval from natural language corpora. At one end of the spectrum, we have ad hoc approaches such as those that compute term co-occurrence frequencies and proximity-based term-term dependencies. And, at the other end of the spectrum, we have more principled approaches, such as those based on Markov Random Fields (MRF) [12], that are based on modeling the term-term dependencies by graphs whose arcs capture the inter-term relationships in the queries vis-à-vis the same in the documents.

With regard to the investigation of such approaches in retrieval from software libraries, in a recent contribution [13], we presented an ad hoc approach that demonstrated how the inter-term proximities can be used to reformulate the queries in order to improve the quality of retrievals. The query reformulation (QR) in that approach takes place through a two-step process that is carried out without any additional input from the user. In the first step, the top-ranked retrievals for the user-supplied query are analyzed for the detection of the terms that are proximal to the query terms as supplied by the user. The proximal terms thus discovered are used to enrich the original query in the second step. We showed that the retrieval for the query reformulated in this way is significantly more precise than it is for the original query.

We now show that even more significant improvement in retrieval precision can be obtained by using a more principled alternative to ad hoc approaches. In particular, we will show that when an MRF is used to model the ordering and the positional dependencies between the query terms vis-à-vis the documents, we end up with a framework that not only yields a higher retrieval precision with the simplest of the ordering and proximity constraints but that can also be generalized to the investigation of more general such constraints. In the MRF-based approach, certain subsets of the terms in a bug report are used for scoring the software artifacts while taking into account term-term proximity and order. This approach exploits the fact that the software artifacts that contain the query terms in the same order and/or in similar proximities as in the query itself are more likely to be relevant to a given query.

While, as demonstrated by our results, MRF is a powerful approach to the modeling of query-document relationships; to fully exploit its potential, it must be used in conjunction with what we refer to as *query conditioning* (QC). The goal of QC is to recognize the fact that a bug report constitutes a highly structured query whose various parts consist, as we mentioned previously, of a textual narrative, a stack trace, etc [11]. These different parts are disparate in the sense that the inter-term relationships do not carry the same weight in them. For example, the proximity of the terms used in the stack trace portions of a bug report carries far more weight than in the textual narrative. Therefore, the ordering and proximity constraints are likely to be far more discriminative in those portions of bug report that, by their very nature, are far more structured. To further underscore the importance of these structured portions of the bug reports, past studies have shown that stack traces are one of the most valuable sources of information for pinpointing the location of the bugs [14, 15].

We have named the overall code retrieval engine that includes both QC and MRF modeling as Terrier+. The reason for the '+' suffix in the name is that ours is an enhancement of the popular open-source research IR engine called Terrier.<sup>1</sup> The functionality we have added to create Terrier+ is highly modular, in the sense that any of the enhancements can be turned on and off merely by clicking buttons in the GUI in order to measure the retrieval effectiveness of the selected

<sup>1</sup><http://terrier.org>

enhancements. This allows for a convenient assessment of the retrieval power that can be attributed to each enhancement.

Note that whereas Terrier+ applies MRF modeling to all queries, QC becomes an important factor only when structured elements are present in the queries. In general, detecting the structured elements in bug reports is a difficult task as they may have different formats and they are usually surrounded by other types of textual data [11]. It is also not uncommon for these constructs to undergo unexpected format changes, such as those caused by accidental line breaks, when they are copied into a bug report. In order to overcome these challenges, Terrier+ employs several regular expressions to detect and extract these structured elements from bug reports.

We experimentally validate the proposed bug localization framework on three large software libraries: AspectJ, Eclipse 3.1 for full bug reports, and Google Chrome for experimenting with bug report titles only. We show that MRF modeling of the queries and the QC step (whenever the queries lend themselves to such conditioning) significantly improve the accuracy with which the bugs can be localized. In order to investigate the effect of the length of queries on the precision with which the bugs are localized, we carried out retrievals with just the bug report titles and with the bug reports taken in their entirety. Whereas MRF modeling resulted in improved precision in bug localization even for short queries consisting of just the bug report titles, the improvements were even more significant when the bug reports in their entirety were subject to MRF modeling and the QC step. Our experimental results also include comparison with the other state of the art IR-based approaches to bug localization. We demonstrate that, on the average, the MRF and QC-based framework outperform all these other approaches.

This paper is organized as follows. In the next section, we start by stating the research questions addressed in the experimental validation of the proposed code retrieval framework. Seeing these questions at the outset will give the reader a better sense of the scope of our work, especially in relation to the previous related contributions by us and by others. We then present the proposed MRF modeling along with QC. We evaluate our retrieval framework in Section 4. Section 5 discusses possible threats to the validity of the experiments we have conducted to establish the power of the proposed approach. The relevant work is presented in Section 6. Finally, we conclude in Section 7.

## 2. RESEARCH QUESTIONS

Our empirical evaluation on large open-source software projects shows that our proposed retrieval framework leads to significant improvements in automatic bug localization accuracy. As mentioned in the previous section, our retrieval engine has been packaged as an enhancement to the popular research IR retrieval engine known as Terrier, and we refer to our enhancement as Terrier+.

We compare the performance of Terrier+ with the other IR approaches developed for the retrieval of the buggy source files in response to bug reports. Included in these other IR-based approaches is the Spatial Code Proximity (SCP)-based QR algorithm in which a given short query is enriched with additional informative terms drawn from the highest ranked retrieval results with respect to the original query for an improved retrieval accuracy [13]. Another important class of IR tools developed for automatic bug localization leverages the past development efforts. Such IR tools have also been shown to improve the bug localization accuracy significantly [8, 9].

For the evaluation of Terrier+, we conducted extensive validation tests with the following research questions in mind:

- RQ1** *Does including code proximity and order improve the retrieval accuracy of bug localization. If so, to what extent?*
- RQ2** *Is QC effective for improving the query representation vis-à-vis the source code?*
- RQ3** *Does including stack traces in the bug reports improve the accuracy of bug localization?*
- RQ4** *How does the MRF-based retrieval framework compare with the Query Reformulation (QR) based retrieval frameworks for bug localization?*

**RQ5** *How does the MRF-based retrieval framework compare with the other bug localization frameworks that leverage the past development history?*

### 3. THE PROPOSED APPROACH

Directly or indirectly, all traditional IR-based approaches to bug localization amount to comparing the first-order distribution of terms in a query vis-à-vis the documents. This applies as much to simple approaches based on VSM and Unigrams [8, 10] as it does to approaches based on Latent Dirichlet Allocation (LDA) [4, 10, 16] that use hidden variables for injecting additional degrees of freedom for comparing the queries with the documents. All of these approaches miss out on the information contained in the inter-term relationships in the queries and in the documents.

While we do now know how to extend the BOW approaches of the sort mentioned previously with proximity-based reformulation of a query for improving the quality of retrievals [13], what we need are theoretically well-grounded approaches that can be generalized to the incorporation of arbitrary inter-term relationships between the terms of a query vis-à-vis the documents. The goal of this section is to address this need by using the notion of MRF.

In the subsections to follow, we first review how the inter-term dependencies are modeled with MRF. We then mention three different specializations of MRF modeling that appear to be particularly appropriate for our needs. Subsequently, in order to exploit MRF modeling to the maximum, we present our QC method to extract from a query those portions that are particularly suited to modeling by MRF. Note that, as we will demonstrate later, MRF improves retrievals even in the absence of QC. However, by giving greater weight to the inter-term relationships in those portions of a query that QC has identified as being highly structured, we improve the quality of retrievals much further.

#### 3.1. Markov Random Fields

Over the years, researchers in the machine learning community have devoted much energy to the investigation of methods for the probabilistic modeling of arbitrary dependencies among a collection of variables. The methods that have been developed are all based on graphs. The nodes of such graphs represent the variables, and the arcs the pairwise dependencies between the variables. The graphs may either be directed, as in Bayesian Belief Networks [17], or undirected, as in the networks derived from MRF [12, 17]. In both these methods, the set of variables that any given variable directly depends on is determined by the node connectivity patterns. In a Bayesian Belief Network, the probability distribution at a node  $q$  is conditioned on only those nodes that are at the tail ends of the arcs incident on  $q$ , taking the *causality* into account. In a Markov network, on the other hand, the probability distribution at a node  $q$  depends on the nodes that are immediate neighbors of  $q$  without considering any directionality. In the context of retrieval from natural language corpora, the work of Metzler and Croft [12] has shown that Markov networks are particularly appropriate for the modeling of inter-term dependencies vis-à-vis the documents.

In general, given a graph  $G$  whose arcs express pairwise dependencies between the variables, MRF modeling of the probabilistic dependencies among a collection  $A$  of variables is based on the assumption that the joint distribution over all the variables can be expressed as a product of nonnegative potential functions over the cliques in the graph:

$$P(A) = \frac{1}{Z} \prod_{k=1}^K \varphi(C_k), \quad (1)$$

where  $\{C_1, C_2, \dots, C_K\}$  represents the set of all cliques in the graph  $G$ , and  $\varphi(C_k)$  a nonnegative potential function associated with the clique  $C_k$ . In the aforementioned expression,  $Z$  is merely for the purpose of normalization because we want the sum of  $P(A)$  over all possible values that can be taken by the variables in  $A$  to add up to unity.

Our end goal with MRF is to rank the files in the code base according to the probability of a file  $f$  in the software library to be relevant to a given query  $Q$ . We denote this probability by  $P(f|Q)$  [8]. Using the definition of the conditional probability, we can write

$$P(f|Q) = \frac{P(Q,f)}{P(Q)}. \quad (2)$$

As we are only interested in ranking the files and considering that the denominator in Equation (2) does not depend on the files, the denominator can be ignored. Hence,  $P(f|Q) \stackrel{rank}{=} P(Q,f)$ . In order to separate out the roles played by the variables that stand for the query terms (because we are interested in the inter-term dependencies in the queries) vis-à-vis the contents of a source file  $f$ , as suggested by Metzler and Croft [12], we will use the following variation of the general form expressed in Equation (1) to compute this joint probability:

$$P(Q,f) = \frac{1}{Z} \prod_{k=1}^K \varphi(C_k) \stackrel{rank}{=} \sum_{k=1}^K \log(\varphi(C_k)), \quad (3)$$

where  $Q$  stands for a query that is assumed to consist of the terms  $q_1, q_2, \dots, q_{|Q|}$  and  $f$  a file in the software library. The nodes of the graph  $G$  in this case consist of the query terms, with one node for each term.  $G$  also contains a node that is reserved for the file  $f$  whose relevancy to the query is in question. As before, we assume that this graph contains the cliques  $\{C_1, C_2, \dots, C_K\}$ . As shown in the formula, for computational ease, it is traditional to express the potential  $\psi(C_k)$  through its logarithmic form, which is through  $\psi(C_k) = \log(\varphi(C_k))$ .

The fact that a fundamental property of any Markov network is that the probability distribution at any node  $q$  is a function of only the nodes that are directly connected to  $q$  may now be expressed as

$$P(q_i|f, q_{j \neq i} \in Q) = P(q_i|f, q_j \in \text{neig}(q_i)), \quad (4)$$

where  $\text{neig}(q_i)$  denotes the terms whose nodes are directly connected to the node for  $q_i$ . As observed in [12], this fact allows arbitrary inter-term relationships to be encoded through appropriate arc connections among the nodes that represent the query terms in the graph  $G$ . At one end of the spectrum, we can assume that the query terms are all independent of one another by the absence of any arcs between them. This assumption, known as the usual BOW assumption in IR, is referred to as the *full independence* (FI) model. And at the other end of the spectrum, we may assume a fully connected graph in which the probability distribution at each node representing a query term depends on all the other query terms (besides being dependent on the file  $f$ ). This is referred to as the *full dependence* (FD) model. Figures 1a and 1c depict the graph  $G$  for FI and FD assumptions for the case when a query  $Q$  consists of exactly three terms.

What makes MRF modeling particularly elegant is that it gives us a framework to conceptualize any number of other ‘intermediate’ forms of dependencies between the two extremes of FI and FD. This we can do by simply choosing graphs  $G$  with different connectivity patterns. Whereas FI is based on the

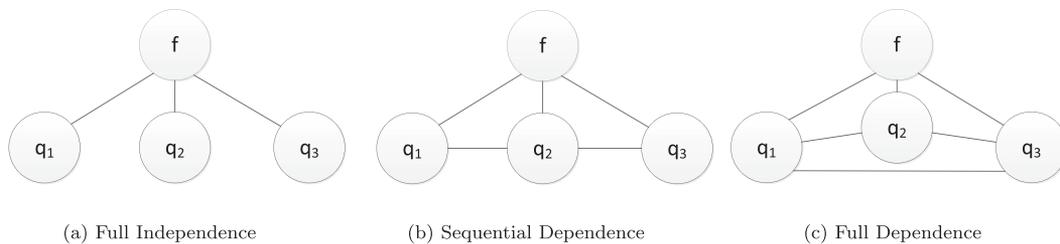


Figure 1. This figure illustrates the notion of a Markov network for capturing the inter-term dependencies in a query vis-à-vis a document. The case depicted is for a query consisting of just three terms. The FI variant shown in (a) corresponds to the Smoothed Unigram Model with retrievals based on Query Likelihoods.

absence of any inter-term arcs in  $G$  and FD on there being an arc between each query term and every other term, we may now think of more specialized dependencies such as the one depicted in Figure 1b. This dependency model, referred to as the *sequential dependency* (SD) model in [12], incorporates both order and proximity between a sequence  $(q_1, q_2, \dots, q_{|Q|})$  of query terms.

At this point, the reader may wonder how one would know in advance as to which connectivity pattern to use for the graph  $G$ . The connectivity pattern is obviously induced by the software library itself. Suppose you believe that the phrase ‘interrupt sig handler’ occurs in the files and can be used to discriminate between them, you would want the nodes for the terms ‘interrupt’, ‘sig’, and ‘handler’ to be connected in the manner shown in Figure 1b.<sup>2</sup>

This is because the SD model shown in that figure would only allow for pairwise (and ordered) occurrences of the words ‘interrupt’ and ‘sig’, on the one hand, and of the words ‘sig’ and ‘handler’, on the other, to be considered when comparing a query with the files. The frequencies with which these terms appear individually in the files may also carry discriminatory power. The importance of a file to a query would be determined by the relative frequencies with which the terms occur individually and the relative frequencies with which the ordered pairs appear in the file vis-à-vis the relative frequencies of the same in the query. In contrast to the case depicted in Figure 1b, should it happen that the queries and the relevant files contain the three terms ‘interrupt’, ‘sig’, and ‘handler’ in all possible orders, you would want to use the FD assumption depicted in Figure 1c. In this case, the number of times these terms occur together within a window of a certain size for each possible order would carry discriminatory power for choosing the files relevant to a query. *If a certain order of appearance for the three terms is prominent in the query, the files that show high counts for the same order of appearance would carry greater weight.* Finally, should it happen that the three terms occur in the relevant files without the order of their appearance being significant in any sense, the BOW assumption implied by the FI model in Figure 1a would be sufficient.

We are particularly interested in the graph connectivity induced by the notion of Spatial Code Proximity (SCP) [13]. SCP consists of first associating a positional index with each term in a query and in the documents as shown in Figure 2. Our goal is to translate the values of the positional indexes into graph models based on the FI, SD, and FD assumptions. In the next three subsections, we will present formulas that show how these models can be derived from SCP-based indexes.

*3.1.1. Full independence.* As already stated, the FI assumption reduces an MRF model to the usual BOW model that has now been extensively investigated for automatic bug localization [8–10]. As should be clear from the graph representation of this model depicted in Figure 1a for the case of a query with exactly three terms, FI modeling involves only 2-node cliques. Therefore, under MRF modeling, the probability of a query given a file is simply computed by summing over the 2-node cliques:  $P_{FI}(f|Q) \stackrel{\text{rank}}{=} \sum_{i=1}^{|Q|} \psi_{FI}(q_i, f)$ . The choice of the potential function, obviously critical in computing this probability, should be in accord with the fact that MRF modeling under the FI assumption amounts to BOW modeling. Therefore, a good choice is to make the potential  $\psi_{FI}(q_i, f)$  proportional to the frequency of the query term  $q_i$  in the file  $f$ . Because the zero probability associated with a query term  $q_i$  that does not appear in a file  $f$  can create problems when estimating the relevance of  $f$  to the query, it is common to add what is referred to as a smoothing increment to the term frequencies. A powerful smoothing approach, known as Dirichlet smoothing, consists of adding to the potential the frequency of the term in the entire corpus [8]. Shown below is a formula for the potential  $\psi_{FI}(q_i, f)$  that includes Dirichlet smoothing:

<sup>2</sup>Since modeling how the terms are *ordered* in a document is central to the SD model, it might seem strange to the reader that the arcs between the nodes that represent the terms are undirected in the graph in that figure. The lack of arrows on the arcs might cause a reader to think that the SD model would not be able to distinguish between the sequences  $(q_1, q_2, q_3)$  and  $(q_3, q_2, q_1)$  of the terms. As we will show later, the term ordering constraints are taken care of by the formulas used for calculating the potentials and the probabilities associated with the different ordering of the terms in the documents vis-à-vis their orderings in the queries. We should also mention that our not using arrows in our graphical depiction of term-term ordering is in keeping with the notational conventions that have emerged over the last three decades in the community of researchers who work on making inferences from network models of multi-variable probabilistic information. These conventions have come into existence in order to distinguish between Markov Random Fields (MRF) and Bayesian Belief Networks.

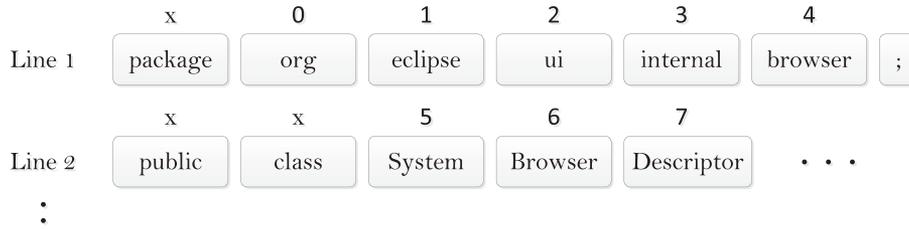


Figure 2. An illustration of indexing the positions of the terms in the Eclipse class SystemBrowserDescriptor.java. The 'x' symbol indicates the stop words that are dropped from the index.

$$\psi_{FI}(q_i, f) = \lambda_{FI} \log \left( \frac{tf(q_i, f) + \mu P(q_i | C)}{|f| + \mu} \right), \quad (5)$$

where  $C$  denotes the entire collection of files in the software library,  $P(q_i | C)$  the probability of the term  $q_i$  in  $C$ ,  $tf(q_i, f)$  the term frequency of  $q_i$  in file  $f$ ,  $|f|$  the length of the file, and  $\mu$  the Dirichlet smoothing parameter. The coefficient  $\lambda_{FI}$  has no impact on the file rankings with this model. However, we keep it in the formulation as we will use it later in SD and FD modeling.<sup>3</sup>

The probability expression shown above for the relevance of a query term to a file is exactly the same as it appears in the widely used BOW model known as *SUM* whose usefulness in automatic bug localization has been demonstrated in [8–10]. We will use the retrieval results obtained with FI as the baseline in order to determine the extent of improvements one can obtain with the other two models, SD and FD.

**3.1.2. Sequential dependence.** The SD model takes the order and the proximity of the terms into account in such a way that the probability law for a query term  $q_i$  given a file  $f$  obeys  $P(q_i | \{q_1, \dots, q_{i-1}, q_{i+1}, \dots, q_{|Q|}\}) = P(q_i | f, q_{i-1}, q_{i+1})$ .

To see how a software library can be processed to induce the SD model, note from the example shown in Figure 1b that we now have 3-node cliques in addition to the 2-node cliques of the FI model. Therefore, we must now count the frequencies with which pairs of terms occur together, with one following the other (without necessarily being adjacent) in a specific order, in addition to counting the frequencies for the terms occurring singly as in the FI model [18]. Again incorporating Dirichlet smoothing for the same reasons as in the FI model, we employ the following potential function for the 3-node cliques corresponding to a file  $f$  and two consecutive query terms  $q_{i-1}$  and  $q_i$ :

$$\psi_{SD}(q_{i-1}, q_i, f) = \lambda_{SD} \log \left( \frac{tf_W(q_{i-1}q_i, f) + \mu P(q_{i-1}q_i | C)}{|f| + \mu} \right), \quad (6)$$

where  $C$  is the collection of all the files in the software library,  $tf_W(q_{i-1}q_i, f)$  the number of times the terms  $q_{i-1}$  and  $q_i$  appear in the same *order* in file  $f$  as in the query within a window length of  $W \geq 2$ . For  $W > 2$ , the terms do not have to be adjacent in the file and a window may also contain other query terms in between. The smoothing increment  $P(q_{i-1}q_i | C)$  is the probability associated with the pair  $(q_{i-1}q_i)$  in the entire software library. To the potential function shown above, we must now add the potential function for 2-node cliques the reader has already seen for the FI model:

$$P_{SD}(f | Q) \stackrel{rank}{=} \sum_{i=2}^{|Q|} \psi_{SD}(q_{i-1}, q_i, f) + \sum_{i=1}^{|Q|} \psi_{FI}(q_i, f). \quad (7)$$

As the reader would expect, the ranking of the files with the combined potential function shown above depends on the relative values of the model parameters  $\lambda_{FI}$  and  $\lambda_{SD}$ . However, the overall scaling of these parameters is inconsequential on account of the unit summation constraint on the

<sup>3</sup>The problems caused by zero-frequency terms in documents and queries have also been dealt with in other related areas, such as traceability link recovery, etc., in software engineering [45] [46].

probabilities. We therefore set  $\lambda_{FI} + \lambda_{SD} = 1$ . We can think of  $\lambda_{SD}$  as an interpolation or a mixture parameter that controls the relative importance of the 3-node cliques vis-à-vis the 2-node cliques.

If the reader is still confused as to how the SD model can distinguish between, say, the sequence (q1, q2, q3) and the sequence (q3, q2, q1) of terms, let us start with the simplest case in which we want to measure the relevancy of a file to a query on the basis of the frequency with which the pair (q1, q2) shows up, with the two terms appearing in that specific order, in the files and in the queries. When we use Equations (6) and (7) in this case for a probabilistic assessment of the relevancy of a file to a bug report, what those equations say is that we must measure the frequency with which the terms (q1,q2) occur together, in that order, in both the files and the bug reports. What that implies is that, in the context of an SD model, we will only measure the frequencies associated with the appearance of (q1,q2) and not with the appearance of (q2, q1). Now let's consider the case where the search window include three terms (q1, q2, q3), and with the SD model, we want to enforce the ordering constraint that q1 appear before q2, and that q2 appear before q3. Using Equations (6) and (7), we measure the frequencies associated with the pairwise appearance of the terms (q1,q2), with the two terms showing up in that order, and with the pairwise appearance of (q2,q3), again, with the two terms showing up in that specific order. Therefore, our frequency of occurrence measurements would ignore the situations where q2 may come before q1 and q3 may come before q2.

*3.1.3. Full dependence.* As demonstrated previously by Figure 1c, the FD assumption implies a fully connected graph  $G$  whose nodes correspond to the individual query terms, with one node being reserved for the file  $f$  under consideration. The graph being fully connected allows for a file  $f$  to be considered relevant to a query regardless of the order in which the query terms occur in the file. (Compare this to the SD case where, for a file  $f$  to be considered relevant to a query, it must contain the query terms in the same order as in the query). Therefore, the FD assumption provides a more flexible matching mechanism for retrievals.

The price to be paid for the generality achieved by FD is the combinatorics of matching all possible ordering of the query terms with the contents of a file. To keep this combinatorial explosion under control, following [18], we again limit ourselves to just 2-node and 3-node cliques. While this may sound the same as for the SD assumption, note that the 3-node cliques are now allowed for a pair of query terms for both orderings of the terms. Therefore, for any two terms  $q_i$  and  $q_j$  of the query, the potential function takes the following form for the 3-node cliques:

$$\psi_{FD}(q_i, q_j, f) = \lambda_{FD} \log \left( \frac{tf_w(q_i q_j, f) + \mu P(q_i q_j | C)}{|f| + \mu} \right), \tag{8}$$

where  $\lambda_{FD}$  again works as a mixture parameter in a manner similar to  $\lambda_{SD}$ ; that is,  $\lambda_{FI} + \lambda_{FD} = 1$ . As before,  $\mu$  is the smoothing parameter, and  $tf_w(q_i q_j, f)$  the frequency for the pair  $q_i q_j$  in  $f$ . Summing over the cliques, we obtain the ranking score of a file by

$$P_{FD}(f|Q) \stackrel{rank}{=} \sum_{i=1}^{|Q|} \sum_{j=1, j \neq i}^{|Q|} \psi_{FD}(q_i, q_j, f) + \sum_{i=1}^{|Q|} \psi_{FI}(q_i, f). \tag{9}$$

*3.1.4. A motivating example.* We will now use a simple example to compare the retrieval effectiveness of the three models: FI, SD, and FD. For reasons of space limitations, we will limit our example to the simplest of the bug reports, one that only contains a one-line text narrative, which corresponds to the title of the bug report.

The bug 98995<sup>4</sup> filed for Eclipse v3.1 has a title that reads ‘monitor memory dialog needs to accept empty expression’. The target source files that were eventually modified to fix this bug are:

1. org.eclipse.debug.ui/.../ui/views/memory/MonitorMemoryBlockDialog.java

<sup>4</sup>[https://bugs.eclipse.org/bugs/show\\_bug.cgi?id=98995](https://bugs.eclipse.org/bugs/show_bug.cgi?id=98995)

Table I. Retrieval accuracies for the Bug 98995 with the three different models.

Method	2-node cliques	3-node cliques	Ranks	AP
FD	7	42	1-3-2	1.0000
SD	7	6	2-3-4	0.6389
FI (BOW)	7	0	6-5-10	0.2778

FI represents the BOW model with the Smoothed Unigram Model.

FI: full independence; BOW: bag-of-words; FD: full dependence; SD: sequential dependence; AP: average precision.

2. org.eclipse.debug.ui/.../ui/views/memory/AddMemoryBlockAction.java
3. org.eclipse.debug.ui/.../ui/DebugUIMessages.java.

After removing the stop word ‘to’ from the title, the final query consists of the following seven terms: ‘monitor memory dialog needs accept empty expression’. Table I presents the retrieval accuracies obtained for this query, along with the number of cliques utilized for each dependency assumption. In the table, the column ‘Ranks’ gives the ranks of the three relevant files in the retrievals. AP is the resulting average precision.

Investigating the ranked lists returned for the three models, we see that FI ranks several irrelevant files above the relevant ones. One such file is ASTFlattener.java. Although this file does not contain any of the terms ‘monitor’, ‘memory’, and ‘dialog’, it is retrieved at the top rank by this model because, as a file related to parsing, it contains the terms ‘accept’, ‘empty’, and ‘expression’ with very high frequencies. Clearly, the FI model misses the word order in the query. Note that the FI model compares only the first-order probability distributions for the term counts.

In comparison with FI, SD is able to retrieve the relevant files at higher ranks, as shown in Table I. The improvement obtained with SD is a consequence of the discriminations achieved by requiring that the query terms, when they appear together in a source file, do so in a specific order. The creation of the 3-node cliques with the query terms is illustrated in Figure 3. A 3-node clique is formed by the two words depicted together with an under-bracket and the node corresponding to a file. Because the relevant files contain these term blocks in close proximity with high frequencies; with this model, they receive higher ranking scores in comparison to the irrelevant files.<sup>5</sup>

Despite the improvements, SD still ranks one irrelevant file, ASTRewriteFlattener.java, above all the relevant ones. This file also does not contain any of the terms ‘monitor’, ‘memory’, and ‘dialog’. However, it contains in close proximity the term pairs for two of the 3-node cliques: ‘accept empty’ and ‘empty expression’. The file manages to receive a high ranking score with these term pairs in addition to the abstract syntax tree related terms.

In this example, FD does better than the other two models because the cliques it incorporates also allow for ordering constraints on pairs of terms that are separated by another term. So if a file has many instances of, say, the pair ‘monitor dialog’, which in the sequence of the seven terms shown earlier is separated by the term ‘memory’, that would contribute to the file becoming more relevant to the query. In general, though, whether or not FD would score higher than SD would depend on how often the differently ordered versions of the same set of terms appear in the files vis-à-vis the bug reports.

### 3.2. Query conditioning

When a bug report contains highly structured components, such as a stack trace and/or a source code patch [11], such information can be crucial to locating the files relevant to the bug [15]. Being highly structured, these components must first be identified as such and subsequently processed to yield the terms that can then be used to form a query for IR-based retrieval. The processing steps needed for

<sup>5</sup>As for the ‘mechanics’ of how the SD model is able to distinguish between the term sequence ‘monitor memory dialog’ and the term sequence ‘memory monitor dialog’, we take the reader back to the explanation at the end of Section 3.1.2. Using Equations (6) and (7), we measure the frequencies associated with the pairwise appearance of the terms in ‘monitor memory’, with the two terms showing up in that order, and with the pairwise appearance of the terms in ‘memory dialog’, again, with the two terms showing up in that specific order. Therefore, our frequency of occurrence measurements ignore the situations where ‘memory’ may come before ‘monitor’ and ‘dialog’ may come before ‘memory’.



Figure 3. Clique creation for the Bug 98995. The figure shows the first four query term blocks for the 3-node cliques utilized by the sequential dependency modeling.

that purpose are different for the two different types of components we consider: stack traces and source code patches. We refer to the collection of these steps as QC. QC is carried out with a set of regular expressions that, while custom designed for the different types of structured components encountered, are sufficiently flexible to accommodate small variations in the structures.<sup>6</sup>

As already stated, our retrieval framework has been packaged as an enhancement to the popular research IR retrieval engine Terrier, and we refer to our enhancement as Terrier+. With regard to the flow of processing related to QC in Terrier+, it uses regular expressions to first identify the patches and the stack traces in a given bug report if any of these elements is available in the report. Then, it processes them separately to sift out the significant source code identifiers to be used in the retrievals. The final query is composed only from the terms extracted from the stack traces and the patches if one or both of these components are available. If these structured components are not available, Terrier+ makes do with the entire bug report as it is and feeds it into the MRF framework.

*3.2.1. Stack traces.* When Terrier+ detects a stack trace in a bug report, it automatically extracts the most likely locations of the bug by identifying the methods in the trace. Since the call sequence in a stack trace starts with the most recent method call,<sup>7</sup> we extract only the topmost  $T$  methods while discarding the rest of the trace; the methods further down in the trace are not likely to contain any relevant terms, and they may introduce noise into the retrieval process. Figure 4 illustrates the stack trace that was included in the report for Bug 77190 filed for Eclipse.<sup>8</sup> The bug caused the *EmptyStackException* to be thrown by the code in *PushFieldVariable.java* and was subsequently fixed in a revised version of this code. The figure highlights the extracted portion of the stack trace that is used in forming the final query. Note that we only extract the methods that are present in the code base to which the bug report applies. That is, we skip the methods from the libraries belonging to the Java platform itself, as illustrated in the figure. During the experiments, we empirically set  $T=3$ , as this setting resulted in the best retrieval accuracies on the average. As we show in our experimental evaluation, this filtering approach increases the precision of the retrievals significantly.

*3.2.2. Patches.* Source code patches are included in a bug report when a developer wishes to also contribute a possible (and perhaps partial) fix to the bug. When contributed by an experienced developer, these components of a bug report can be directly used for pinpointing the files relevant to a bug.

A patch for a given bug is usually created with the *Unified Format* to indicate the differences between the original and the modified versions of a file in a single construct. With this format, the textual content of the patch contains the lines that would be removed or added in addition to the contextual lines that would remain unchanged in the file after the patch is applied. For term extraction from the patches, Terrier+ does not use the lines that would be added after the suggested patches are applied to the files as those lines are not yet present in the code base.

Obviously, the files mentioned by a developer in a patch may not correspond to the actual location of the bug. And, there may be additional files in the code base that may require modifications in the final fix for the bug. While the importance of information in such source code patches cannot be

<sup>6</sup>Our QC only takes into account the stack traces and source code patches when they can be identified in a bug report. Note that a bug report may also contain additional source code snippets that are not meant to be patches [5]. QC treats any additional such code on par with the main textual part of the report.

<sup>7</sup>We do realize that for some languages the methods in a stack trace are in the opposite order. That is, the most recent method call appears as the last entry in the trace. The logic of identifying the methods most relevant to a bug would obviously need to be reversed for such languages.

<sup>8</sup>[https://bugs.eclipse.org/bugs/show\\_bug.cgi?id=77190](https://bugs.eclipse.org/bugs/show_bug.cgi?id=77190)

```

java.util.EmptyStackException
  at java.lang.Throwable.<init>(Throwable.java)
  at java.util.Stack.peek(Stack.java)
  at java.util.Stack.pop(Stack.java)
  at org.eclipse.jdt.internal.debug.eval.ast.engine.Interpreter.pop(Interpreter.java:89)
  at org.eclipse.jdt.internal.debug.eval.ast.instructions.Instruction.popValue(Instruction.java:111)
  at org.eclipse.jdt.internal.debug.eval.ast.instructions.PushFieldVariable.execute(PushFieldVariable.java:54)
  at org.eclipse.jdt.internal.debug.eval.ast.engine.Interpreter.execute(Interpreter.java:50)
  ⋮
  at org.eclipse.core.internal.jobs.Worker.run(Worker.java:66)

```

Figure 4. The stack trace that was included in the report for Bug 77190 filed for Eclipse. With query conditioning, the trace is first detected in the report with regular expression based processing. Subsequently, the highlighted lines are extracted as the most likely locations of the bug and fed into the Markov Random Field framework.

overstated, it is important to bear in mind that their inclusion in the bug reports is more the exception than the rule. As mentioned in Table III, out of the 4035 bug reports we analyzed for Eclipse v3.1, only eight contained a patch. Along the same lines, out of the 291 bug reports we analyzed for the AspectJ project, only four contained a patch. Nonetheless, considering the importance of the information contained in the patches when they are included in a bug report, Terrier+ takes advantage of that information whenever it can.

Figure 5 illustrates the data flow in the retrieval framework with QC and MRF.

#### 4. EXPERIMENTAL EVALUATION

We evaluate the effect of incorporating term-term dependencies on the retrievals for bug localization for three large software projects, namely Eclipse IDE<sup>9</sup> and AspectJ<sup>10</sup> for experimenting with full bug reports, and Google Chrome<sup>11</sup> for experimenting with just the bug report titles. We evaluate QC on only Eclipse and AspectJ as the bug reports for Chrome do not contain stack traces or patches. We use a set of bug reports, denoted  $B$ , that were filed for these projects and, for the ground truth, the files modified to fix the corresponding bugs as the relevant file sets to be retrieved by the retrieval engine. The relevant file set for a bug report is denoted  $RF$ .

Because the bug tracking databases such as Bugzilla<sup>12</sup> do not usually store the modification histories of the changes made to the files in response to the bug reports, researchers commonly use the commit messages in the repository logs in order to establish the links between the repository modifications and the bug reports [9, 13, 19–21]. The *BUGLinks*<sup>13</sup> [13] and the *iBugs*<sup>14</sup> [19] are the resulting datasets of such approaches that give us the links between the bug reports and the files relevant to the bugs in the repositories for the projects we have used. The *BUGLinks* dataset contains information related to the Eclipse and the Chrome projects, whereas the *iBugs* dataset contains information related to the AspectJ project. For Eclipse and Chrome, we indexed the version 3.1 and 4.0.305.0, respectively, to perform the retrievals, whereas for AspectJ, we indexed the corresponding version for each bug report separately as the bugs in the *iBugs* dataset are filed for different versions of this project. Tables II and III present various statistics drawn from these datasets regarding the three projects used in our evaluation study. In Table II,  $|B|$  denotes the number of bug reports used in querying the code base of each project,  $avg(|RF|)$  the average number of relevant files per bug, and  $avg(|Q_{Title}|)$  the average length of the bug report titles that are used in the retrievals.<sup>15</sup> Additional information for the two cases in which we also used QC for retrieval is presented in Table III. In the table, #Patches

<sup>9</sup>[www.eclipse.org](http://www.eclipse.org)

<sup>10</sup><http://eclipse.org/aspectj/>

<sup>11</sup>[www.google.com/chrome](http://www.google.com/chrome)

<sup>12</sup><https://bugs.eclipse.org/bugs/>

<sup>13</sup><https://engineering.purdue.edu/RVL/Database/BUGLinks/>

<sup>14</sup><https://www.st.cs.uni-saarland.de/ibugs/>

<sup>15</sup>As we describe in Section 4.3, our experimental studies involve two types of experiments: those that are based on just the titles of the bug reports, and those that include the titles and the bug descriptions. Obviously, whereas MRF modeling can be applied to both types of queries, query conditioning (QC) can only be investigated for the latter type.

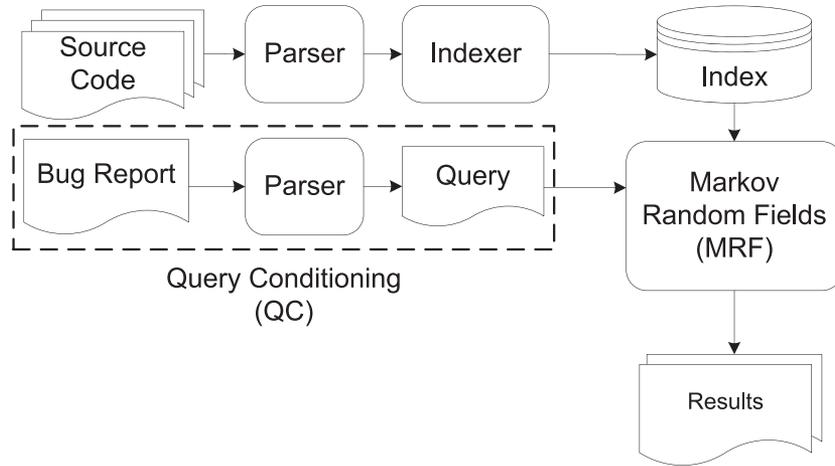


Figure 5. An illustration of the data flow in the proposed retrieval framework.

Table II. Evaluated projects.

Project, description	Language	B	$avg( RF )$	$avg( Q_{Title} )$
AspectJ, an extension to Java programming language	Java	291	3.09	5.78
Eclipse v3.1, integrated development environment	Java	4,035	2.76	5.80
Chrome v4.0 (for titles only), web browser	C/C++	358	3.82	6.21

Table III. Various statistics related to the bug reports used in the experiments for the evaluation of bug localization with Markov Random Field modeling and query conditioning.

Project	#Patches	#ST	#Remaining	$ Q_{Title + Desc} $
AspectJ	4	89	198	56.77
Eclipse v3.1	8	519	3,508	44.11

and #ST show the number of bug reports containing patches and stack traces, respectively, and  $|Q_{Title + Desc}|$  is the average length of the bug reports, including both the title and the description without any filtering, in terms of the number of tokens used in querying the code base. In order to evaluate the accuracy of extracting structural elements with QC, we randomly sampled 1000 bug reports for the Eclipse project and manually investigated them for the presence or the absence of structural elements. Because AspectJ has a relatively smaller set of bug reports, we manually examined all of them to evaluate the accuracy of structural element identification for this project.  $Pr$  and  $Re$  in Table IV show the precision and the recall numbers for detecting the stack traces. The precision and recall for detecting the patches were both 1.0.

In the rest of this section, we will start out with how the source files and the bug reports are tokenized for the extraction of the terms that are used to represent each file and each bug report. Subsequently, we describe the metrics we use in our evaluation study. That is followed by the experimental results demonstrating the power of the MRF-based approach to the modeling of the source code libraries.

Table IV. The accuracy with which stack traces are detected in bug reports prior to MRF based retrieval.

Project	Pr	Re
AspectJ	1.0	0.96
Eclipse v3.1	1.0	0.93

#### 4.1. Preprocessing of the source code files and the bug reports

For the indexing of a particular version of the target code base, we first split the compound terms that are typically formed by using punctuation characters and camel casing. Then we drop the programming language specific terms and a set of standard English stop words. The remaining terms are then stemmed to their common roots using the Porter's stemming algorithm [22]. The position of each term extracted from the files is recorded after these preprocessing steps, as illustrated in Figure 2. These steps constitute the front end to the enhancements that distinguish Terrier+ from Terrier. As mentioned previously, Terrier+ is an extension to the open-source research search engine Terrier [23]. As for the bug reports, they are also subject to the same preprocessing steps.

Subsequent to preprocessing, Terrier+ represents a file by a multidimensional array that can be accessed via its ID. This data structure contains term IDs, the corresponding term frequencies, and the positions of the terms in the file. For each term in a given query, the files that contain the term are accessed via an *inverted index* in which a term is represented by a two dimensional array that stores the file IDs and the frequency of the term in those files.

#### 4.2. Evaluation metrics

We evaluate the retrieval accuracy of Terrier+ using precision and recall based metrics [22]. We have tabulated the bug localization performance using precision at rank  $r$  ( $P@r$ ), recall at rank  $r$  ( $R@r$ ), and mean average precision (MAP) metrics. While  $P@r$  measures the accuracy with which the files are retrieved up to rank  $r$ ,  $R@r$  measures the completeness of the retrievals. The average precision (AP) for a query  $Q \in B$  is given by

$$AP(Q) = \frac{\sum_{r=1}^{RT} P@r \times I(r)}{rel_Q}, \quad (10)$$

where  $I(r)$  is a binary function whose value is 1 when the file at rank  $r$  is a relevant file, and 0 otherwise. The parameter  $RT$  in the summation bound stands for the total number of highest-ranked files that are examined for the calculation of  $AP$  for a given query  $Q$ . The denominator  $rel_Q$  is the total number of relevant files in the collection for  $Q$ .  $AP$  estimates the area under the precision-recall curve for a given query  $Q$ , and therefore, it is suitable for comparing retrieval algorithms [22]. MAP is computed by taking the mean of the APs for all the queries. In addition to these metrics; we also present the number of hits ( $H@r$ ) for the bug reports [24], which gives the number of bugs for which at least one relevant source file is retrieved in the ranked lists above a certain cut-off point  $r$ .

We use MAP for comparing the different retrieval methods as it is the most comprehensive metric that takes into account both the precision and the recall at multiple ranks. In computation of this metric for the results we report in this paper, we set  $RT=100$  in the summation in Equation (10). In order to evaluate whether the improvements obtained with the proposed approaches are significant or not, we use pairwise student's  $t$ -test on the APs for the queries.

#### 4.3. Bug localization experiments

For an in-depth analysis of the retrievals, we divide each bug report into two parts, *Title* and *Description*. We first conduct two sets of experiments using these two parts for each bug report *without QC*: (1) retrievals with MRF modeling using only the titles of the bug reports. The queries used for these retrievals are denoted 'title-only'. And (2) retrieval with MRF modeling using the complete bug reports; that is, including both the titles and the descriptions for the bug reports. The queries used for these retrievals are denoted 'title+desc'. Then, we incorporate QC in the second category of retrievals and analyze the usefulness of including stack traces and patches in queries by comparing the overall retrieval accuracy for the set of bug reports that contain these elements to the remaining set of the bug reports in our query set.

**4.3.1. Parameter sensitivity analysis.** The model parameters that affect the quality of the retrievals in our retrieval framework are: (1) the window length parameter ( $W$ ); (2) the mixture parameters of the

respective dependency models ( $\lambda_{SD}, \lambda_{FD}$ ); and (3) the Dirichlet smoothing parameter ( $\mu$ ). While  $W$  sets the upper bound for the number of intervening terms between the terms of the 3-node cliques,  $\lambda_{SD}$  and  $\lambda_{FD}$  simply adjust the amount of interpolation of the scores obtained with the 2-node cliques with those obtained with the 3-node cliques as explained in Section 3. As the ranking is invariant to a constant scaling in the mixture parameters, we enforce the constraints  $\lambda_{FI} + \lambda_{SD} = 1$  and  $\lambda_{FI} + \lambda_{FD} = 1$  for SD and FD modeling, respectively. In all experiments, we empirically set the Dirichlet smoothing parameter as  $\mu = 4000$ . Note that the retrieval accuracy is not very sensitive to the value of this parameter [8].

Figures 6 and 7 show retrieval accuracies for bug localization in terms of MAP as the window length and the mixture parameters are varied for the ‘title-only’ and the ‘title+desc’ queries. As shown in Figures 6a and 7a, a value of 0.2 works well for the mixture parameters in general. Note that when  $\lambda_{SD} = \lambda_{FD} = 0.0$ , SD and FD use only the 2-node cliques, and hence, they reduce to FI, the SUM. As for the window length, Figures 6b and 7b illustrate the effect of varying this parameter for  $\lambda_{SD} = \lambda_{FD} = 0.2$ . On the average,  $W = 8$  results in high retrieval accuracies for the analyzed projects for both types of queries.

As shown in Figures 6b and 7b, when the window length is set as  $W = 2$ , SD performs better than FD in all experiments across the projects. This is because the terms are required to be adjacent to be matched in the code base when we use this setting, and therefore, the order of the terms becomes more important.

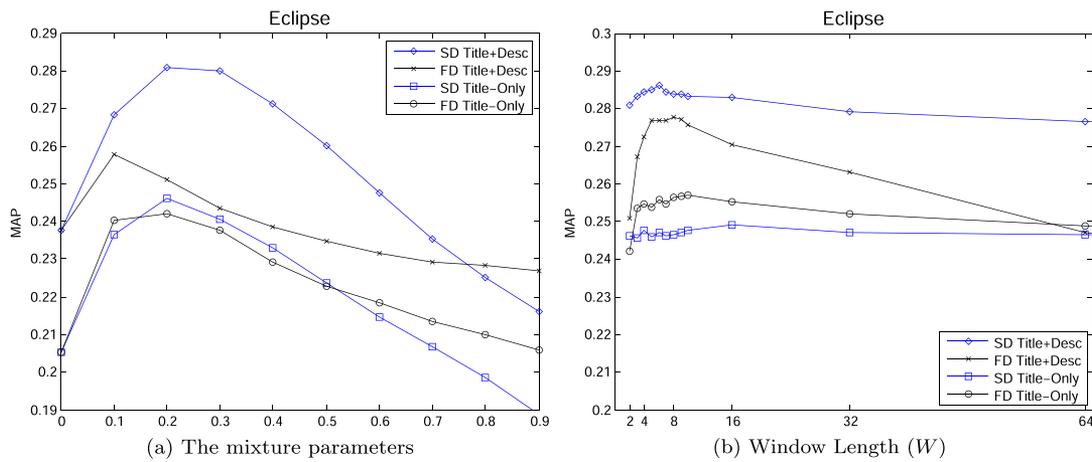


Figure 6. The effects of varying model parameters on mean average precision (MAP) for Eclipse. The figure on the left shows the MAP values as the mixture parameters ( $\lambda_{SD}$  for the sequential dependency assumption and  $\lambda_{FD}$  for the full dependency assumption) are varied while the window length parameter is fixed as  $W = 2$ . The figure on the right shows the MAP values as  $W$  is varied while the mixture parameters are fixed at 0.2.

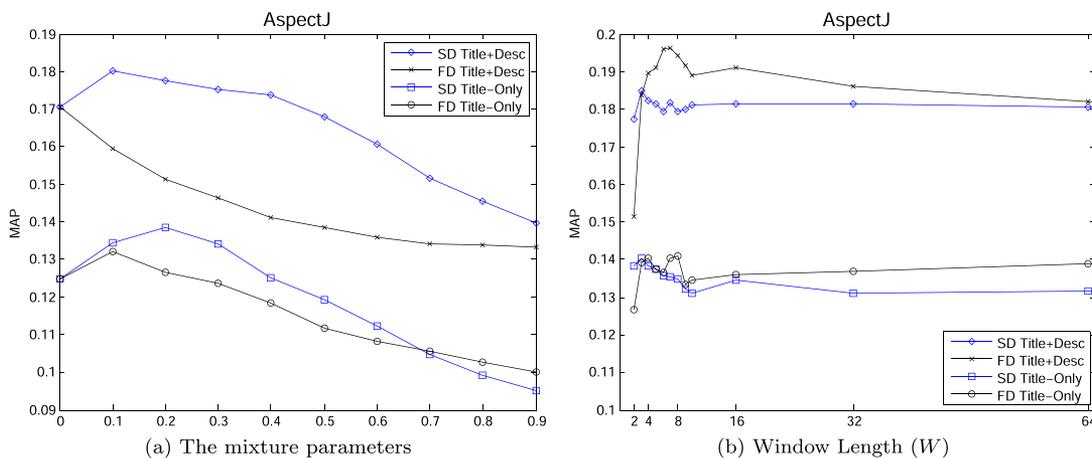


Figure 7. The effects of varying model parameters on mean average precision for AspectJ. The rest of the caption here is the same as for Figure 6.

Overall, SD is less sensitive to the window length parameter, achieving high retrieval accuracies consistently.

*4.3.2. Retrieval results.* In this section, we compare the retrieval performances of the two term-term dependency models, SD and FD, with the FI model. We fix the interpolation parameters as  $\lambda_{SD} = \lambda_{FD} = 0.2$  and the window lengths as  $W=8$  in all the experiments presented in the remainder of this paper.

Table V presents the bug localization accuracies on the evaluated projects for the ‘title-only’ queries and MRF modeling. With these experiments, we explore the retrieval accuracy of Terrier+ for short queries comprised of only a few terms. The last row of the table shows the ‘baseline’ accuracy; this is obtained with the FI assumption, which, as mentioned previously, is the same thing as SUM [8–10]. The highest score in each column is shown in bold. All the dependency-model based improvements reported in this table over the FI model are statistically significant at  $\alpha=0.05$  level. Note that incorporating the term dependencies into the retrievals improves the accuracy of bug localization substantially in terms of the six metrics presented in the table.

Table VI presents the bug localization accuracies for the ‘title+desc’ queries without QC. (The results obtained with QC will be shown in Subsection 4.3.3 that follows.) That is, the entire textual content of the bug reports, without any QC, is used for querying the code base. The improvements obtained with FD and SD over FI are statistically significant at  $\alpha=0.05$  in this table. Note that the retrieval accuracies improve significantly when the description parts of the bug reports are also included for carrying out the retrievals (even in the absence of any QC). While SD and FD perform

Table V. Retrieval accuracy with the ‘title-only’ queries. (The baseline, FI represents the BOW assumption).

Eclipse						
Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.2564 (+24.83%)</b>	<b>0.2198</b>	<b>0.1110</b>	<b>0.3199</b>	<b>0.4070</b>	<b>2,100</b>
SD	0.2466 (+20.06%)	0.2116	0.1069	0.3083	0.3934	2,042
FI (BOW)	0.2054	0.1710	0.0883	0.2556	0.3417	1,805
AspectJ						
Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.1410 (+12.89%)</b>	<b>0.1409</b>	<b>0.0832</b>	<b>0.1794</b>	<b>0.2420</b>	124
SD	0.1348 (+7.93%)	0.1340	0.0790	0.1675	0.2382	<b>125</b>
FI (BOW)	0.1249	0.1375	0.0708	0.1498	0.2079	111

MAP: mean average precision; FD: full dependence; SD: sequential dependence; FI: full independence; BOW: bag-of-words.

Table VI. Retrieval accuracy for the ‘title+desc’ queries.

Eclipse						
Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	0.2778 (+16.87%)	0.2496	0.1201	0.3427	0.4317	2,249
SD	<b>0.2840 (+19.48%)</b>	<b>0.2543</b>	<b>0.1232</b>	<b>0.3530</b>	<b>0.4391</b>	<b>2,268</b>
FI (BOW)	0.2377	0.2020	0.1060	0.3046	0.3859	2,044
AspectJ						
Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.1945 (+14.08%)</b>	<b>0.2131</b>	<b>0.0997</b>	<b>0.2322</b>	0.2996	142
SD	0.1794 (+5.22%)	0.1856	0.0990	0.2203	<b>0.3096</b>	<b>148</b>
FI (BOW)	0.1705	0.1856	0.0880	0.2003	0.2693	126

MAP: mean average precision; FD: full dependence; SD: sequential dependence; FI: full independence; BOW: bag-of-words.

comparably well in these experiments on the Eclipse project, FD outperforms SD on AspectJ on average in terms of MAP.

We are now in a position to answer the first of the five research questions (RQ) formulated in Section 2 of this paper. For convenience, here is the question again:

**RQ1** *Does including code proximity and order improve the retrieval accuracy for bug localization. If so, to what extent?*

Based on the results presented in Tables V and VI, we conclude that incorporating spatial code proximity and order in the retrieval process improves the accuracy of automatic bug localization significantly. On average, for both short and long queries, which may contain stack traces and/or patches, SD and FD modeling consistently enhance the retrieval performance of Terrier+ over FI across the projects. The improvements are up to 24.83% for Eclipse and up to 14.08% for AspectJ in terms of MAP.

**4.3.3. The effect of query conditioning on retrievals.** The retrieval accuracies obtained with QC on the ‘title + desc’ queries are presented in Table VII, where each bug report is first probed for stack traces and patches in order to extract the most useful source code identifiers to be used in bug localization as explained in Section 3.2. The results shown in Table VII help us answer the following research question that was presented in Section 2:

**RQ2** *Is QC effective on improving the query representation vis-à-vis the source code?*

Comparing the results presented in Tables VI and VII, we conclude that QC indeed leads to superior query formulation for source code retrieval. For all three forms of the dependency assumption, we obtain significant improvements with QC in terms of the six evaluation metrics mentioned in the tables.

The main benefits of QC are seen for the bug reports that contain stack traces because patches are included only in a few bug reports. Figure 8 presents the retrieval accuracies of Terrier+ obtained specifically with the bug reports that contain stack traces. The figure shows that, on the average, the accuracy of the retrievals doubles with QC for both projects in terms of MAP, reaching values above the 0.3 threshold.

Figure 8 also demonstrates the effect of the MRF modeling with stack traces. Comparing the results obtained with the two term-term dependency models, we observe that FD and SD outperform FI consistently when QC is used in retrievals with the stack traces. The main reason for these results is that the order and the proximity of the terms in stack traces are extremely important in locating the relevant source files. As we also mentioned in Section 3, the likelihood of a file to be relevant to a query increases when it contains longer phrases from the stack trace with the same order and proximity relationships. Interestingly, when the queries are not processed with QC, the average retrieval accuracy with FD on the Eclipse project is slightly lower than FI. This is clearly due to the

Table VII. Retrieval accuracy for the ‘title + desc’ queries with query conditioning.

Method	Eclipse					
	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.3019 (+17.93%)</b>	0.2696	0.1274	0.3709	<b>0.4640</b>	<b>2,386</b>
SD	0.3014 (+17.74%)	<b>0.2704</b>	<b>0.1290</b>	<b>0.3749</b>	0.4599	2,354
FI (BOW)	0.2560	0.2186	0.1114	0.3263	0.4102	2,147
Method	AspectJ					
	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.2307 (+8.31%)</b>	<b>0.2715</b>	<b>0.1155</b>	<b>0.2703</b>	0.3400	161
SD	0.2263 (+6.24%)	0.2646	0.1148	0.2658	<b>0.3554</b>	<b>167</b>
FI (BOW)	0.2130	0.2440	0.1052	0.2438	0.3215	147

MAP: mean average precision; FD: full dependence; SD: sequential dependence; FI: full independence; BOW: bag-of-words.



Figure 8. The effect of query conditioning on bug localization with bug reports containing stack traces.

noise in the lengthy stack traces that contain many method signatures most of which are irrelevant to the bug. The QC framework effectively removes the irrelevant method signatures from the trace for a better query representation.

4.3.4. *The role of stack traces in automatic bug localization.* Studies have shown that the developers who are in charge of fixing bugs look for stack traces, test cases, and the steps contained therein, in order to reproduce the bugs, these being the most useful structural elements for comprehending the underlying cause of the bugs and fixing them [14, 15]. Among these structural elements, stack traces are very important for the work we report in this paper. They are not only frequently included in bug reports but also a good source of discriminative source code identifiers for automatic bug localization. Figure 9 presents the retrieval accuracies obtained with the bug reports containing different types of structural elements. In the figure, ‘remaining’ denotes the bug reports that do not contain any stack traces or patches.

That sets us up to answer the research question RQ3 that was previously stated in Section 2:

**RQ3** *Does including stack traces in the bug reports improve the accuracy of automatic bug localization?*

As demonstrated in Figure 9, bug reports with patches lead to the highest accuracies as expected. After the patches, stack traces hold the second position in terms of their usefulness in locating the relevant source code. The retrieval results for the remaining bug reports that do not contain any stack traces or patches are the worst in this comparison. Based on these results, we conclude that including stack traces in the bug reports does improve the bug localization accuracy. Note that the retrieval accuracies we obtained with the stack traces are consistently above the 0.3 threshold for the analyzed projects in terms of MAP.

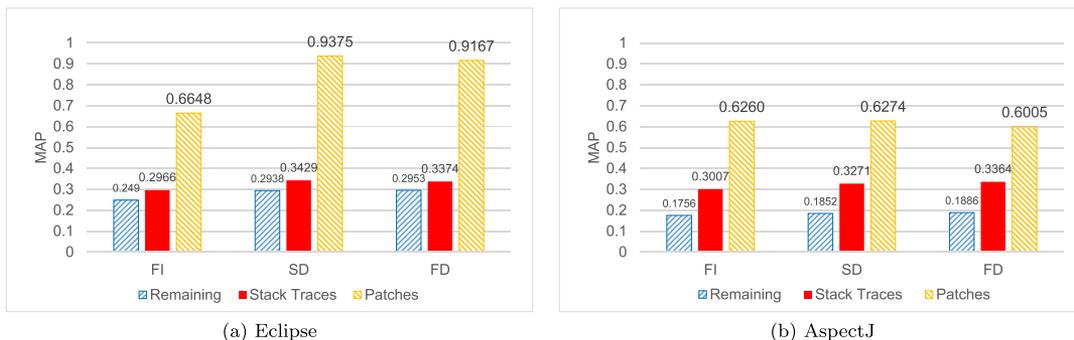


Figure 9. The effect of including structural elements in bug reports on automatic bug localization accuracy. Patches lead to the highest retrieval scores while the bug reports with no stack traces or patches perform the worst in terms of the mean average precision. The number of queries used in these experiments are given in Table III.

#### 4.4. Comparison to automatic query reformulation

In [13], we proposed an automatic QR framework to improve the query representation for bug localization. For experimental evaluation, we used the title of a bug report as an initial query, which was reformulated via pseudo relevance feedback based on the retrieval results obtained with the initial query. The experimental evaluation of the approach showed that the proposed SCP-based QR model in [13] outperformed the state-of-the-art QR models. That admittedly very brief introduction to SCP-based QR leads us to the next question that was posed originally in Section 2:

**RQ4** *How does the MRF-based retrieval framework compare with the QR-based retrieval framework for bug localization?*

Comparing the retrieval accuracies presented in Tables VIII and IX to the retrieval accuracies of the QR models reported in [13], we observe that MRF framework outperforms the SCP-based QR (denoted as SCP-QR in the tables) on the average. For the Chrome project, while the differences between the AP obtained with the respective models are not statistically significant at  $\alpha=0.05$ , the differences in terms of the presented recall metrics are. Additionally, H@10 values obtained with the MRF framework are considerably higher than the values obtained with SCP-QR. For the Eclipse project, both SD and FD perform better than SCP-QR in terms of the metrics shown. The differences are statistically significant at  $\alpha=0.05$ .

#### 4.5. Comparison with bug localization techniques that use prior development history

Another important class of IR approaches to bug localization is based on the prior development history [8, 9]. In [9], Zhou *et al.* proposed BugLocator, a retrieval tool that uses the textual similarities between a given bug report and the prior bug reports to enhance the bug localization accuracy. The main motivation behind BugLocator is that the same files tend to get fixed for similar bug reports during the life-cycle of a software project. Another study that leverages the past development efforts was reported by us in [8]. In that work, we mined the software repositories for the defect and modification likelihoods of the source files in order to estimate a prior probability distribution, which could then be used for more accurate source code retrieval for bug localization. This brief review of this class of approaches to bug localization takes us to the last of the research questions posed in Section 2:

**RQ5** *How does the MRF-based retrieval framework compare with the other bug localization frameworks that leverage the past development history?*

Table VIII. Query Reformulation versus Markov Random Field for the Eclipse project with the ‘title-only’ queries.

Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.2564</b>	<b>0.2198</b>	<b>0.1110</b>	<b>0.3199</b>	<b>0.4070</b>	<b>2,100</b>
SD	0.2466	0.2116	0.1069	0.3083	0.3934	2,042
SCP-QR	0.2296	0.1906	0.1014	0.2853	0.3746	1,915

MAP: mean average precision; FD: full dependence; SD: sequential dependence; SCP-QR: spatial code proximity - query reformulation.

Table IX. Query Reformulation versus Markov Random Field for the Chrome project with the ‘title-only’ queries.

Method	MAP	P@1	P@5	R@5	R@10	H@10
FD	<b>0.1951</b>	<b>0.1844</b>	<b>0.1061</b>	<b>0.2394</b>	<b>0.3159</b>	<b>178</b>
SD	0.1814	0.1760	0.1039	0.2288	0.3137	177
SCP-QR	0.1820	0.1788	0.0933	0.2021	0.2775	151

MAP: mean average precision; FD: full dependence; SD: sequential dependence; SCP-QR: spatial code proximity - query reformulation.

The accuracy of BugLocator was also evaluated on Eclipse v3.1 and iBugs datasets. The evaluation on the Eclipse project was carried out using 3075 bug reports filed for version 3.1 of the software, whereas our experiments with Terrier+ were performed using 4035 bug reports filed for the same version. In order to compare the performance of Terrier+ with that of BugLocator, we repeated the experiments using only the bug reports with which BugLocator was evaluated.

Using the MAP metric, Figure 10 shows the accuracy of the proposed framework along with that of BugLocator. In the figure, we also include the accuracies obtained with the revised Vector Space Model (rVSM) [9] that, according to the authors of BugLocator, yields retrieval results superior to those obtained with the classic Vector Space Model (VSM). As can be seen in the figure, FD+QC and SD+QC outperform BugLocator with MAP values above 0.32 threshold for the Eclipse project. In comparison, BugLocator performs better than SUM and SUM+QC with a MAP value of 0.30, while SUM+QC outperforms rVSM. The performance comparisons with the different models are similar for the AspectJ project.

The other bug localization technique that reports improved retrieval performance with past development history leverages the version histories for a software library [8]. In that study, we showed that the Term Frequency - Inverse Document Frequency (TFIDF) model incorporating defect histories for the software artifacts (TFIDF + Defect History based Prior with Decay (DHbPd)) reaches a MAP value of 0.2258 on the AspectJ project.

Based on these results, we conclude that Terrier+ performs at least as good as these other bug localization techniques *but without having to leverage the past development history for a software library*.

#### 4.6. Comparison with Lobster

In this section, we compare our approach with the Lobster (LOcating Bugs using Stack Traces and tExt Retrieval) framework presented in [25].

Lobster uses a combination of textual similarity and stack trace based similarity to retrieve source code files relevant to bug reports. In order to compute the textual similarity between a source code file and a bug report, the method uses the traditional BOW model (VSM). On the other hand, stack trace based similarity is computed using the distance between the source code file and the file present in the stack trace. Note that this distance is calculated by, first, creating the dependency graph of the library, and then, finding the shortest path (in the dependency graph) between each source code file and the file present in the stack trace. The dependency graph is a directed graph in which each node represents a distinct source code file and each edge from one node to another represents control or data flow. If there is no path between two files in the dependency graph, the distance is assumed to be infinity. For such cases, the stack trace based similarity is set to zero, and the total similarity (which is a linear combination of textual and stack trace based similarities) relies only on the textual similarity.

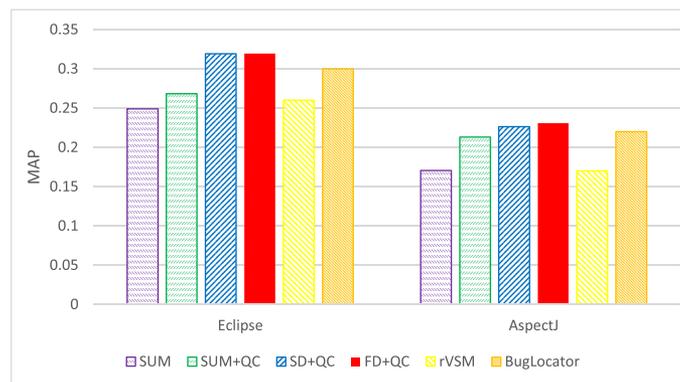


Figure 10. Comparison of the retrieval models for bug localization.

The two important parameters of Lobster are  $\alpha$  and  $\lambda$ . The value of  $\alpha$  adjusts the relative weights of the textual and stack trace based similarities, whereas the value of  $\lambda$  indicates how far a source code file can be from the stack trace file in the dependency graph in order for the two to be considered similar. The values of  $\alpha$  and  $\lambda$  were set to 0.8 and 2, respectively, as recommended by the authors of Lobster.

To compare our MRF approach with Lobster, we installed the software provided by its authors in one of our computers and realized that it did not come with the preprocessing code needed for stop word removal, stemming, and so on. We therefore had to add our own preprocessing code to the downloaded software. Initially, though, we applied the downloaded code to the ZooKeeper *preprocessed data* that the authors have made available at their website. This was done to make sure that our installation of Lobster was correct. Running our install of Lobster on the downloaded preprocessed files for ZooKeeper yielded a MAP value of 0.66, which is roughly the same as 0.63 in the Lobster publication.

What we found interesting was that when we preprocessed the ZooKeeper source code libraries with our own best-practice algorithms, the MAP values obtained using the Lobster software dropped down from 0.66 to 0.48. (Note that we used the same stop word list as provided by the authors at their website [25].) The main conclusion to be drawn from these MAP values is that a significant component of retrieval precision is predicated on how the preprocessing of the source code libraries is carried out. It is also important to note that there is yet no standardized set of preprocessing algorithms for the source files for the kind of research reported here. The MAP values obtained for the installation verification procedure carried out by us are tabulated in the top portion of Table X.

Now we will report on a comparison of Lobster with our MRF-based approach on the AspectJ library. For this comparison, we consider only those queries that contain stack traces. As shown in the lower part of Table X, the MRF-based approach significantly outperforms Lobster. The former yields a MAP of 0.34, as was shown previously in Figure 8b, whereas the same for the latter is 0.05. Note that we use only those queries that contain stack traces for this comparison in order to conform to the evaluation methodology used in the Lobster paper.

## 5. THREATS TO VALIDITY

The threats to the validity of our approach mainly emanate from the scope of the datasets we used in our evaluations and the procedures used to create them. Although our experimental evaluation involves large open source projects commonly used in the evaluation of bug localization approaches, the performance of the proposed framework may vary for other open source or propriety projects. An important step used in the preparation of these datasets is the reconstruction of the links between the bug tracking databases and the corresponding development effort in the software repositories. This reconstruction step is performed using regular expressions to link the bug reports to the repository commits based on the commit messages. Despite the fact that a large number of bug reports are accurately linked to the repository commits, this linking process may occasionally fail when the commit messages in the versioning tools are much too cryptic for regular-expression based matching to succeed [26].

Table X. Lobster versus Markov Random Field for bug localization.

	Method	Preprocessing by	Result generated by	Library	MAP
<b>Installation Verification</b>	Lobster	Moreno <i>et al.</i>	Moreno <i>et al.</i> (in publication)	ZooKeeper	0.63
	Lobster	Moreno <i>et al.</i>	Us	ZooKeeper	0.66
	Lobster	Us	Us	ZooKeeper	0.48
<b>Comparison</b>	MRF	Us	Us	AspectJ	<b>0.34</b>
	Lobster	Us	Us	AspectJ	0.05

MAP: mean average precision.

## 6. RELEVANT WORK

Traditional methods for bug localization rely on the dynamic or the static properties of software [27–32]. Whereas dynamic approaches require that a set of test cases be executed to locate the parts of the program causing the bug, static approaches aim to leverage the static properties of the software such as its function call graphs, dependency relationships between the code segments, etc. The main problem with the static approaches is that they tend to return too many false positives [33]. Although dynamic approaches tend to be more accurate than static methods, designing an exhaustive set of test cases to reveal defective behaviors is very difficult and expensive. The bug localization approach we presented in this paper does not call for a program to be executed. Besides, it is lightweight and inexpensive.

Concept location, feature/concern location, and bug localization are closely related problems in software engineering. Early work on using text retrieval methods for concept location includes the work by Marcus *et al.* [5]. They used Latent Semantic Indexing (LSI) to retrieve the software artifacts in response to short queries. The retrievals are performed in the lower dimensional LSI space, which assigns greater importance to the terms that frequently co-occur in the source files. This framework can also be used to expand a given initial query that consists of just a single query term. In [34], Poshyvanyk *et al.* extended this approach to include formal concept analysis. They showed that the irrelevant search results returned by the LSI model can be reduced with formal concept analysis. Hybrid methods that combine dynamic analysis with IR have also been proposed in this area [2, 35].

In [36], Hill *et al.* leveraged source code identifiers to automatically extract the phrases relevant to a given initial query. These phrases were then used to either find the relevant program elements or to manually reformulate the query for superior feature/concern localization. In [37], Hill *et al.* investigated the effect of the position of a query term on the accuracy of the search results. The main idea in this study is that the location of a query term in the method signatures and in the method bodies determines its importance in the search process. Along the same lines, Ripon *et al.* [38] have improved the precision of the basic BOW-based bug localization by giving greater weight to ‘structural terms’. This contribution is based on the premise that there is much structural information in source code, and it ought to be exploited in a retrieval algorithm for bug localization. The structure that Ripon *et al.* have used is reflected in the terms that are used for classes, methods, and variables. An abstract syntax tree representation of a source code file (which can be obtained by parsing the file) immediately yields the terms that are used as the names for classes, methods, and variables. These terms can subsequently be weighted differentially to reflect their importance in bug localization. Suppose a bug report explicitly mentions a class, the differential weighting given to the terms that stand for classes are likely to result in higher-precision retrieval of the relevant files. As to how this work compares with ours, note that the work of Ripon *et al.* does not directly model the term-term dependencies. Yes, indirectly, by giving differential weighting to the terms used for, say, the classes, the work by Ripon *et al.* introduces an element of positional specificity into the retrieval process – but that is only with respect to a select category of terms (such as those used for classes). The overall retrieval framework in Ripon *et al.* remains that of BOW. Our work, on the other hand, goes beyond BOW and uses MRF modeling to capture proximity relationships between arbitrary terms in the files. Yet another approach to capturing proximity and order relationships in text files – but only at the character level – is through  $n$ -grams-based modeling of a software library [39]. Our work captures such relationships at a higher level of abstraction – the word level.

We should also mention a recent study by Hill *et al.* [40] that investigates the effect of incorporating term-term proximity and order on feature location. This work shows that MRF-based approaches to leveraging positional proximity leads to more consistent retrieval accuracies in comparison to NLP-based approaches [40].

In [41], Gay *et al.* used explicit relevance feedback for QR for the purpose of concept location. This framework requires developers to engage in an iterative query/answer session with the search engine. At each iteration, the developer is expected to judge the relevance of the returned results vis-à-vis the current query. Based on these judgments, the query is reformulated with the Rocchio formula [42] and

resubmitted to obtain the next round of retrieval results. This process is repeated until the target file is located or the developer gives up.

In [43], Haiduc *et al.* introduced Refocus, an automatic QR tool for text retrieval in software engineering. Refocus automatically reformulates a given query by choosing the best QR technique, which is determined by training a decision tree on a separate query set and their retrieval results. After training, based on the statistics of the given query, the decision tree recommends an automatic QR technique that is expected to perform the best among the others.

Recently, several studies have investigated the IR algorithms for the retrieval of software artifacts for bug localization. In a comparative study, Rao and Kak evaluated a number of generic and composite IR models to localize the files that should be fixed to resolve bugs [10]. The main result that came out of this study was that simpler models, such as the VSM or the SUM, performed better than the more sophisticated models such as LDA. Another similar comparative study is by Lukins *et al.* [4]. Their results show LDA performing at least as well as Latent Semantic Analysis (LSA). In a related contribution, Nguyen *et al.* [16] also proposed an LDA-based approach to narrow down the search space for improving bug localization accuracy. In [7], Ashok *et al.* have shown how the relationship graphs can be used to retrieve source files and prior bugs in response to what they refer to as ‘fat queries’ that consist of structured and unstructured data. And, finally, Thomas *et al.* [44] have shown that parameter tuning has a significant effect on retrieval accuracy. The study also showed that combining multiple models can further enhance the retrieval accuracy of even the best performing models.

## 7. CONCLUSIONS

We presented a theoretically sound IR framework for automatic bug localization that takes into account spatial code proximity and term ordering relationships in a code base for improved retrieval accuracy. The proposed retrieval framework benefits from a fuzzy matching mechanism between the term blocks of the queries and the source code. At the heart of the approach is the concept of MRF that captures both the positional and the order attributes of the terms in source files. Our experimental validation involving large open-source software projects and over 4000 bugs has conclusively established that the retrieval performance improves significantly when both the proximity and ordering relationships between the terms are taken into account in matching bug reports with files.

Our experimental evaluation also demonstrates that, in conjunction with the MRF model, the proposed QC approach effectively exploits the different types of structural information that is frequently included in bug reports. We showed that the structural elements in bug reports – particularly stack traces – contain vital information that is not well represented vis-à-vis the source code with the widely used BOW assumption.

Overall, on the basis of retrieval accuracies, it is clear that one can obtain significantly higher bug localization accuracies when MRF modeling and QR are included in a retrieval framework compared with all other state-of-the-art approaches, even including those that take prior development history into account.

## ACKNOWLEDGEMENTS

This work was supported by Infosys.

## REFERENCES

1. Gethers M, Dit B, Kagdi H, Poshyvanyk D. Integrated impact analysis for managing software changes. *Software Engineering (ICSE)*, 2012 34th International Conference on. IEEE, 2012; 430–440.
2. Liu D, Marcus A, Poshyvanyk D, Rajlich V. Feature location via information retrieval based filtering of a single scenario execution trace. *Proceedings of the Twentysecond IEEE/ACM International Conference on Automated Software Engineering*. ACM, 2007; 234–243.
3. Lucia AD, Fasano F, Oliveto R, Tortora G. Recovering traceability links in software artifact management systems using information retrieval methods. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 2007; **16**(4):13.

4. Lukins S, Kraft N, Etkorn L. Source code retrieval for bug localization using Latent Dirichlet Allocation. *Reverse Engineering*, 2008. WCRE'08. 15th Working Conference on. IEEE, 2008; 155–164.
5. Marcus A, Sergeev A, Rajlich V, Maletic J. An information retrieval approach to concept location in source code. *Reverse Engineering*, 2004. Proceedings. 11th Working Conference on. 2004; 214–223. DOI 10.1109/WCRE.2004.10
6. Poshyvanyk D, Marcus A, Ferenc R, Gyimóthy T. Using information retrieval based coupling measures for impact analysis. *Empirical Software Engineering* 2009; **14**(1):5–32.
7. Ashok B, Joy J, Liang H, Rajamani SK, Srinivasa G, Vangala V. DebugAdvisor: A recommender system for debugging. Proceedings of the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering. ACM, 2009; 373–382.
8. Sisman B, Kak A. Incorporating version histories in information retrieval based bug localization. *Mining Software Repositories (MSR)*, 2012 9th IEEE Working Conference on. IEEE, 2012; 50–59.
9. Zhou J, Zhang H, Lo D. Where should the bugs be fixed? More accurate information retrieval-based bug localization based on bug reports. *Software Engineering*, 2012 34th International Conference on. IEEE, 2012; 14–24.
10. Rao S, Kak AC. Retrieval from software libraries for bug localization: A comparative study of generic and composite text models. *Proceeding of the 8th Working Conference on Mining Software Repositories*. 2011; 43–52.
11. Bettenburg N, Premraj R, Zimmermann T, Kim S. Extracting structural information from bug reports. *Proceedings of the 2008 International Working Conference on Mining software repositories*. ACM, 2008; 27–30.
12. Metzler D, Croft W. A Markov random field model for term dependencies. *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2005; 472–479.
13. Sisman B, Kak AC. Assisting code search with automatic query reformulation for bug localization. *Proceedings of the 10th Working Conference on Mining Software Repositories, MSR '13*. IEEE Press, Piscataway, NJ, USA, 2013; 309–318. URL <http://dl.acm.org/citation.cfm?id=2487085.2487145>
14. Bettenburg N, Just S, Schröter A, Weiss C, Premraj R, Zimmermann T. What makes a good bug report? *Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of Software Engineering*. ACM, 2008; 308–318.
15. Schroter A, Bettenburg N, Premraj R. Do stack traces help developers fix bugs? *Mining Software Repositories (MSR)*, 2010 7th IEEE Working Conference on. IEEE, 2010; 118–121.
16. Nguyen AT, Nguyen TT, Al-Kofahi J, Nguyen HV, Nguyen T. A topic-based approach for narrowing the search space of buggy files from a bug report. *Proceedings of 26th IEEE/ACM International Conference on Automated Software Engineering (ASE'11)*. 2011; 263–272.
17. Kollar D, Friedman N. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press, 2009.
18. Peng J, Macdonald C, He B, Plachouras V, Ounis I. Incorporating term dependency in the DFR framework. *Proceedings of the 30th annual international ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2007; 843–844.
19. Dallmeier V, Zimmermann T. Extraction of bug localization benchmarks from history. *Proceedings of the twenty-second IEEE/ACM international conference on Automated software engineering*. ACM, 2007; 433–436.
20. Nguyen AT, Nguyen TT, Nguyen HA, Nguyen TN. Multi-layered approach for recovering links between bug reports and fixes. *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12*. ACM, New York, NY, USA, 2012; 63:1–63:11. DOI 10.1145/2393596.2393671.
21. Wu R, Zhang H, Kim S, Cheung S. Relink: recovering links between bugs and changes. *SIGSOFT FSE*. 2011; 15–25.
22. Manning C, Raghavan P, Schütze H. *Introduction to Information Retrieval*, vol. **1**. Cambridge University Press: Cambridge, 2008.
23. Macdonald C, He B, Plachouras V, Ounis I. University of Glasgow at TREC 2005. Experiments in terabyte and enterprise tracks with Terrier. *Proceedings of TREC 2005*. 2005.
24. Nguyen AT, Nguyen TT, Al-Kofahi J, Nguyen HV, Nguyen TN. A topic-based approach for narrowing the search space of buggy files from a bug report. *26th International Conference on Automated Software Engineering (ASE'11)*. IEEE, 2011; 263–272.
25. Moreno L, Treadway J, Marcus A, Shen W. On the use of stack traces to improve text retrieval-based bug localization. *IEEE International Conference on Software Maintenance and Evolution*. 2014.
26. Bachmann A, Bird C, Rahman F, Devanbu P, Bernstein A. The missing links: Bugs and bug-fix commits. *Proceedings of the eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering*. ACM, 2010; 97–106.
27. Dallmeier V, Lindig C, Zeller A. Lightweight bug localization with AMPLE. *Proceedings of the 6th International Symposium on Automated Analysis-driven Debugging*. ACM, 2005; 99–104.
28. Hovemeyer D, Pugh W. Finding bugs is easy. *ACM SIGPLAN Notices* 2004; **39**(12):92–106.
29. Jones JA, Harrold MJ, Stasko J. Visualization of test information to assist fault localization. *Proceedings of the 24th International Conference on Software Engineering, ICSE '02*. ACM, New York, NY, USA, 2002; 467–477. DOI 10.1145/581339.581397.
30. Liblit B, Naik M, Zheng AX, Aiken A, Jordan MI. Scalable statistical bug isolation. *ACM SIGPLAN Notices* 2005; **40**(6):15–26.
31. Liu C, Fei L, Yan X, Han J, Midkiff SP. Statistical debugging: A hypothesis testing-based approach. *IEEE Transactions on Software Engineering* 2006; **32**(10):831–848.

32. Robillard MP. Topology analysis of software dependencies. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 2008; **17**(4):18.
33. Dit B, Revelle M, Gethers M, Poshyvanyk D. Feature location in source code: A taxonomy and survey. *Journal of Software: Evolution and Process* 2013; **25**(1):53–95. DOI:10.1002/smr.567.
34. Poshyvanyk D, Gethers M, Marcus A. Concept location using formal concept analysis and information retrieval. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 2012; **21**(4):23.
35. Poshyvanyk D, Guéhéneuc YG, Marcus A, Antoniol G, Rajlich V. Combining probabilistic ranking and latent semantic indexing for feature identification. 14th IEEE International Conference on Program Comprehension. IEEE, 2006; 137–148.
36. Hill E, Pollock L, Vijay-Shanker K. Automatically capturing source code context of NL-queries for software maintenance and reuse. 31st International Conference on Software Engineering, 2009. ICSE 2009, 2009; 232–242. DOI 10.1109/ICSE.2009.5070524
37. Hill E, Pollock L, Vijay-Shanker K. Improving source code search with natural language phrasal representations of method signatures. Proceedings of the 2011 26th IEEE/ACM International Conference on Automated Software Engineering. IEEE Computer Society, 2011; 524–527.
38. Saha R, Lease M, Khurshid S, Perry D. Improving bug localization using structured information retrieval. 2013 IEEE/ACM 28th International Conference on Automated Software Engineering (ASE). 2013; 345–355. DOI 10.1109/ASE.2013.6693093
39. Lal S, Sureka A. A static technique for fault localization using character n-gram based information retrieval model. Proceedings 5th India Software Engineering Conference. 2012; 109–118.
40. Hill E, Sisman B, Kak AC. On the use of positional proximity in IR-based feature location. CSMR-WCRE. 2014; 318–322.
41. Gay G, Haiduc S, Marcus A, Menzies T. On the use of relevance feedback in IR-based concept location. Software Maintenance, 2009. ICSM 2009. IEEE International Conference on. IEEE, 2009; 351–360.
42. Rocchio J. Relevance feedback in information retrieval. 1971.
43. Haiduc S, Bavota G, Marcus A, Oliveto R, De Lucia A, Menzies T. Automatic query reformulations for text retrieval in software engineering. Proceedings of the 2013 International Conference on Software Engineering, ICSE '13. IEEE Press, Piscataway, NJ, USA, 2013; 842–851. URL <http://dl.acm.org/citation.cfm?id=2486788.2486898>
44. Thomas S, Nagappan M, Blostein D, Hassan A. The impact of classifier configuration and classifier combination on bug localization. *IEEE Transactions on Software Engineering* 2013; **39**(10):1427–1443. DOI:10.1109/TSE.2013.27.
45. Antoniol G, Canfora G, Casazza G, Lucia AD, Merlo E. Recovering traceability links between code and documentation. *IEEE Transactions on Software Engineering* 2002; **28**:970–983.
46. Linares-Vasquez M, White M, Bernal-Cardenas C, Moran K, Poshyvanyk D. Mining Android app usages for generating actionable GUI-based execution scenarios. Proceedings 12th Working Conf. on Mining Software Repositories. 2015.