# Using Consensual Biterms from Text Structures of Requirements and Code to Improve IR-Based Traceability Recovery

Hui Gao
ghalexcs@gmail.com
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

Hongyu Kuang
khy@nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

Kexin Sun
mf20320130@smail.nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

Xiaoxing Ma
xxm@nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

Alexander Egyed
alexander.egyed@jku.at
Institute for Software Systems
Engineering, Johannes Kepler
University
Linz, Austria

Patrick Mäder
patrick.maeder@tu-ilmenau.de
Fakultät für Informatik und
Automatisierung, Technische
Universität Ilmenau
Ilmenau, Germany

Guoping Rong
ronggp@nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

Dong Shao
dongshao@nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

He Zhang
hezhang@nju.edu.cn
State Key Lab of Novel Software
Technology, Nanjing University
Nanjing, China

## ABSTRACT

Traceability approves trace links among software artifacts based on whether two artifacts are related by system functionalities. The traces are valuable for software development, but are difficult to obtain manually. To cope with the costly and fallible manual recovery, automated approaches are proposed to recover traces through textual similarities among software artifacts, such as those based on Information Retrieval (IR). However, the low quality & quantity of artifact texts negatively impact the calculated IR values, thus greatly hindering the performance of IR-based approaches. In this study, we propose to extract co-occurred word pairs from the text structures of both requirements and code (i.e., consensual biterms) to improve IR-based traceability recovery. We first collect a set of biterms based on the part-of-speech of requirement texts, and then filter them through the code texts. We then use these consensual biterms to both enrich the input corpus for IR techniques and enhance the calculations of IR values. A nine-system-based evaluation shows that in general, when solely used to enhance IR techniques, our approach can outperform pure IR-based approaches and another baseline by 21.9% & 21.8% in AP, and 9.3% & 7.2% in MAP, respectively. Moreover, when used to collaborate with another enhancing strategy from different perspectives, it can outperform this baseline by 5.9% in AP and 4.8% in MAP.

## CCS CONCEPTS

• **Software and its engineering** → **Software development process management**.

## KEYWORDS

traceability recovery, text structures, biterm, information retrieval

## 1 INTRODUCTION AND MOTIVATION

Software traceability is "the ability to interrelate any uniquely identifiable software engineering artifact to any other, maintain required links over time, and use the resulting network to answer questions of both the software product and its development process"[1]. For example, these *traces* can reveal how stakeholders' expectations of system functionalities (i.e., requirements [14]) are actually implemented and executed during the running of the system (i.e., code [9]). Recent work has shown that, when correctly recovered, requirements-to-code traces can help developers to perform software maintenance tasks faster and more correctly[39], and the software quality is also highly relevant to the completeness of these traces[61]. In practice, traceability is not only mandated in certain regulations to demonstrate that a system is safely running [40, 48], but also increasingly used to help ensure the security of a

---

[1]http://www.CoEST.org

system [47, 50]. Unfortunately, the manual recovery of traceability is labor-intensive and error-prone [22, 58] due to the semantic gap between artifacts of different concept levels [8] (such as requirements and code), and a large number of traces to recover.

To both improve the overall accuracy of traceability and reduce manual efforts, a growing body of automated approaches is proposed to recover or maintain traces based on mainly information retrieval (IR) and machine learning (ML) techniques. To bridge the semantic gap between different artifacts, these approaches typically use the calculated textual similarities of artifacts to represent the likelihood of the actual traces. Particularly, when recovering traces from scratch, i.e., no prior identified traces are available to form the training set of ML techniques, the semi-automated, IR-based approaches become the mainstream in research [15]. Instead of verifying all possible traces between any two artifacts, users of these approaches can verify candidate trace links along the automatically generated candidate lists sorted by textual similarities (in descending order) that are calculated through IR models, such as Vector Space Model (VSM) [6], Latent Semantic Indexing (LSI) [43], and the probabilistic Jensen and Shannon model (JS) [1].

Unfortunately, different software artifacts (such as requirements and code) often use different terms to denote the same concept because they are at different concept levels [8]. This situation, which is known as the *vocabulary mismatch problem*, greatly hinders the performance of IR techniques on traceability recovery. Naturally, researchers proposed to use advanced lexical analyses to either enhance the similarity calculation through IR models [24, 51], or enrich the artifact texts as the input of IR models [16, 20]. However, the improvement brought by these approaches is limited to the unsatisfying quantity and quality of artifact texts, especially for real-world software projects. Therefore, researchers further proposed enhancing strategies that are built upon other unique perspectives of IR-based traceability recovery, such as incorporating with execution tracing [21, 57] by running the systems, combining with the analyses on code structural information (e.g., call and data dependencies between code elements) [35, 44], and utilizing user's verification on candidate links [18, 23, 28, 52, 53] due to the semi-automated nature of IR-based traceability recovery. But yet the performance of these enhancing strategies are still highly relevant to the quality of initially calculated IR similarities.

Meanwhile, a different body of work aims to reduce the "noise" from the texts of different artifact texts by focusing on the writing of software documents and code. De Lucia et al. [19] proposed a smoothing filter to remove the duplicate information in the same kind of artifacts (e.g., requirements or code) that does not carry the semantics of the artifact. Furthermore, Capobianco et al. [12] proposed that software artifacts are typically written by a "technical language" because the users of this language, i.e., the developers from the same software project, work in a particular area and have a common interest [30, 32]. In this case, the nouns contained in the artifact texts are more likely to carry the semantics of the artifact and thus are more valuable to improve IR-based traceability recovery. Recently, Ali et al. [2] reported that not only nouns, but also the other major Part-of-Speech (POS) categories (including verbs, adjectives, adverbs, and pronouns) are important for IR-based traceability recovery. They further prune generated IR candidate links by finding whether the requirement and code from each given

candidate link share at least one verb, to improve IR-based traceability recovery. Although these approaches alleviate the vocabulary mismatch problem by reducing the "noise" in artifact texts, they do not enhance, or even weaken the underlying artifact semantics (e.g., only consider nouns for IR techniques), thus preventing them from bringing in more improvement.

Unlike these approaches, in this paper, we propose to use *consensual biterms* that are extracted from both the requirements and the code, to first enrich the artifact texts as the input for IR-based traceability recovery, and then use these extracted biterms to further improve the ranking of generated IR candidate lists. In particular, a *biterm* is an unordered term-pair co-occurring over a term sequence within a text. It is first proposed to address the sparse data problem for retrieving documents through IR techniques [62], and then used in topic modeling over short texts [13] to address the same problem. Specifically, we first extract an initial set of biterms from the requirement texts based on the grammatical relationships between two terms in each sentence by conducting the Stanford typed dependencies parsing [42]. We then extract another set of biterms from code identifiers and comments by using a sliding window. Finally, we keep the biterms from the intersection of these two sets and name them as **consensual biterms** because we argue that these biterms indicate at least part of the same system functionalities that are characterized in different software artifacts at different concept level, like requirements and code elements reach a "consensus" over these biterms. We then argue that extracting consensual biterms from common software projects is viable because as we discussed, requirements are written through a technical language that is though less formal than programming language, but still more rigorous than common natural language, and code is implemented through a well-defined programming language and the naming of its identifiers usually follows suggested conventions. For example, the Java Language Specification [2] recommends that class names should be nouns or noun phrases, and method names should be verbs or verb phrases (C Language also has the similar naming conventions [3]). Therefore, we propose that *the consensual biterms extracted from the texts of both requirements and code can help to bridge the semantic gap between requirements and code, and thus they are very important to improve IR-based traceability recovery.* We will use the following example (adapted from the widely researched iTrust system [4]) for further demonstration.

iTrust is a medicare system and one of its requirements (UC35) is to report adverse events on prescription drugs or immunization reactions through emails (though the email is never sent because iTrust is designed for lecturing testing courses only). Accordingly, the class EmailUtil is implemented to send those emails and thus it is traced to UC35. However, as depicted in Figure 1, UC35 only has one sub-flow S1 to describe the use of email. Thus, the textual similarity between UC35 and EmailUtil is inevitably small when calculated through IR models. Furthermore, the term "send" and its past participle "sent" are often filtered by the stop word list of IR-based approaches for traceability recovery because they appear too frequently in typical software artifacts and is considered as part of the "noise" in artifact texts. Unfortunately, without this term, the

---

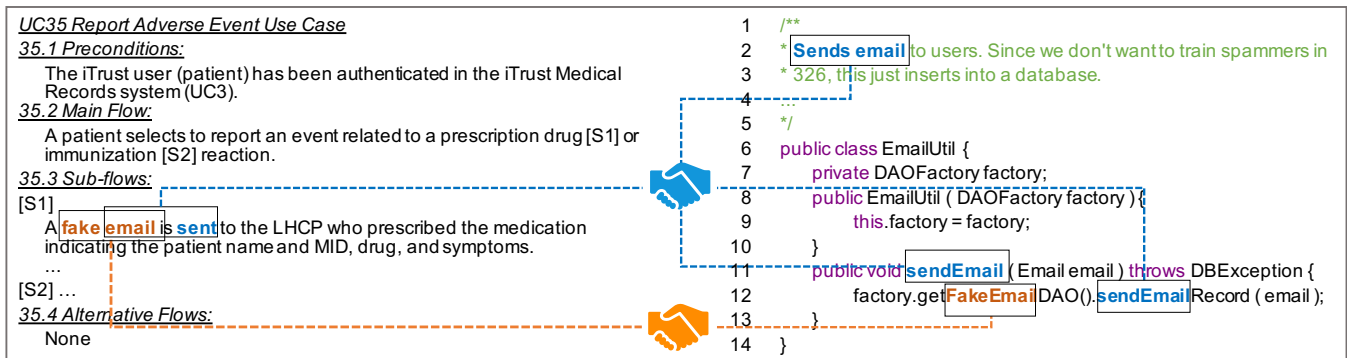| | |
|---|---|
| *UC35 Report Adverse Event Use Case* | 1  /** |
| *35.1 Preconditions:* | 2  * **Sends email** to users. Since we don't want to train spammers in |
| The iTrust user (patient) has been authenticated in the iTrust Medical Records system (UC3). | 3  * 326, this just inserts into a database. |
| *35.2 Main Flow:* | 4  ... |
| A patient selects to report an event related to a prescription drug [S1] or immunization [S2] reaction. | 5  */ |
| | 6  public class EmailUtil { |
| *35.3 Sub-flows:* | 7      private DAOFactory factory; |
| [S1] | 8      public EmailUtil ( DAOFactory factory ) { |
| A **fake email** is **sent** to the LHCP who prescribed the medication indicating the patient name and MID, drug, and symptoms. | 9          this.factory = factory; |
| ... | 10     } |
| [S2] ... | 11     public void **sendEmail** ( Email email ) throws DBException { |
| *35.4 Alternative Flows:* | 12         factory.get**FakeEmail**DAO().**sendEmail**Record ( email ); |
| None | 13     } |
| | 14 } |

**Figure 1: Motivating example adapted from the iTrust system**

system functionality "send email" is also omitted by IR models as well, making the rank of EmailUtil very low in the IR candidate list for UC35. In contrary, we use the Stanford typed dependencies parsing to extract two biterms from the sub-flow S1, i.e., (send, email) and (fake, email), based on two grammatical relationships: *noun subject* and *adjectival modifier*, respectively. These two biterms can also be extracted from the method identifiers and the class comment of EmailUtil, and they become the consensual biterms. We argue that these two consensual biterms can narrow down the semantic gap between UC35 and EmailUtil by carrying the shared system functionality: "send fake email to report adverse event", thus helping to improve IR-based traceability recovery.

Based on the consensual biterm, we further propose an enhancing strategy for IR-based, requirements-to-code traceability recovery. We first extracts consensual biterms from the texts of requirements and code, and then use them to enrich the corpus with other extracted terms for the following similarity calculations through IR models. Because consensual biterms carry the shared system functionalities between requirements and code, we further use them to adjust the generated IR value for a give candidate link by counting the occurrence of a consensual biterm globally (i.e., between a requirement and a code file) and locally (i.e., among different sections of a requirement, e.g., "title" and "sub-flow" for use cases, while "summary" and "description" for issues). Eventually, all IR candidate lists are re-ranked based on both the enriched corpus with consensual biterms and the adjustments of IR candidate links based on the occurrences of consensual biterms in code and the text structures of requirements. The evaluation of our proposed approach is based on nine real-world systems, and also involves the three mainstream IR models (i.e., VSM, LSI, and JS). Our evaluation first showed that our approach can statistically outperform the pure IR-based approach and the state-of-the-art IR-based approach that utilizes the POS tags for enhancement [2]. Our evaluation then showed that our approach, which is proposed by analyzing the writing of artifact texts, can also collaborate with another enhancing strategy [23] based on two different perspectives (i.e., combining user feedback with code dependency analysis), to further improve IR-based traceability recovery.

The contribution of this paper is extracting consensual biterms from the texts of requirements and code to both enrich the input corpora for IR models and improve the rankings of IR candidate

lists. This work mainly targets functional requirements and focuses on tracing the requirements to code classes particularly. We name our approach as **TAROT** (**T**race**A**bility **R**ecovery by c**O**sensual bi**T**erms). This work contains two novel features: (1) we extract consensual biterms from requirements and code based on grammatical relationships (by using the Stanford CoreNLP toolkit) and naming conventions, respectively; (2) we improve the overall accuracy of IR-based traceability recovery through extracted consensual biterms to both enrich the input corpora for IR techniques and re-rank the IR candidate lists.

## 2  BACKGROUND AND RELATED WORK

This section discusses the background and related work on typical approaches for traceability recovery, enhancing strategies for IR-based approaches, and the use of biterms in lexical analyses.

**IR-based traceability recovery:** Information-Retrieval (IR) - based approaches are perhaps now the most representative ones in automated traceability recovery [15]. These approaches use IR techniques to compute the textual similarities between two different software artifacts (e.g., requirements and code), and suggest the IR values as the probability of whether a pair of source-target artifacts is a relevant trace. IR techniques compute the textual similarity based on the occurrence of terms in the text of artifacts. If two artifacts share a larger number of terms, they will be more textually similar, and thus are more likely to be linked by a relevant trace. In general, typical IR-based approaches use three steps to compute textual similarities [17]: (1) they build a corpus with terms extracted from different software artifacts (the terms from code are usually extracted from identifiers and comments only) after pre-processing such as stemming and stop words removal; (2) these approaches use the term-by-document matrix to represent the corpus where each artifact are organized as a document in the matrix where its cell values are usually weighted by the tf-idf weighting scheme [7] to distinguish the importance for each term in the document based on its the occurring frequency; (3) these approaches compute the similarity among artifacts (represented as entries in the term-by-document matrix) through different IR models. The Vector Space Model (VSM) [6] treats each document as a vector of terms. In this model, the similarity of two artifacts is typically calculated as the cosine of the angle between the two vectors in the space of terms occurred in all artifacts. The Latent

Semantic Indexing (LSI) [43] is a VSM that applies Singular Value Decomposition (SVD) [7] to the term-by-document matrix to construct an LSI subspace. The size of the subspace is a manually tuned threshold called the k value. Because the LSI subspace captures the most significant factors through SVD, LSI is expected to implicitly handle the most frequently co-occurring terms and thus to alleviate the vocabulary mismatch problem for IR-based traceability recovery. Meanwhile, the Jenson-Shannon model (JS) [1] treats each document as a probability distribution containing the probability of a term occurring in each document. JS calculates the similarity among documents through the Jenson-Shannon divergence, i.e., the "distance" between two given probability distributions. To evaluate whether TAROT is generally beneficial to IR-based approaches, our evaluation involves all three discussed IR models.

**ML-based traceability recovery:** A growing body of work proposes to take advantage of Machine Learning techniques for traceability recovery. These ML-based approaches can work well especially when part of the trace links have been previous identified, such as completing the missing links between issues and code commits in software repository [36, 60, 63], or maintaining recovered traces through machine learning classification [46]. It is worth-while noting that due to the data imbalance problem between known traces and no-traces, re-balancing steps (e.g., undersampling no-traces or oversampling traces) are necessary for these ML-based approaches. These approaches often use calculated IR values as the features of their classification, while our approach is expected to improve the quality of IR values, thus being likely to improve these approaches as well. Researchers also proposed deep-learning-based approaches, such as the one proposed by Guo et al. [25] that uses the RNN network to generate links between subsystem requirements and design definitions, and T-BERT [37] that uses different BERT architectures to recover missing links between issues and commits. The deep learning (DL) network can capture the implicit connections from term sequences in artifact texts, while our approach extracts biterms explicitly from term sequences in software texts, and code identifiers and comments. However, the DL-based approaches still requires correctly-labeled training and development sets to work. To address this issue, Mills et al. [45] introduce active learning to substantially reduce the amount of required training data for classification approaches while maintaining similar performance. Moran et al. [47] use a tailored hierarchical Bayesian Network to infer candidate traces through transitive relationships among different groups of software artifacts, and their approach Comet only needs a small amount of user feedback for better inference. Although we extract consensual biterms from the texts of both requirements and code, the extraction itself does not need any trace links to be previously identified. Thus, our approach, which is an enhancement for IR-based approaches, can still work when the user has to recover trace links from scratch.

**Enhancement for IR-based approaches:** To address the vocabulary mismatch problem that greatly hinders the performance of IR-based traceability recovery, researchers have proposed many enhancing strategies from different perspectives, such as introducing enhanced lexical analyses on the text of artifacts [16, 19, 24], incorporating with execution tracing [21, 57], or combining with the analyses on different kinds of code dependencies [35, 44]. Meanwhile, because the semi-automatic nature of IR-based traceability
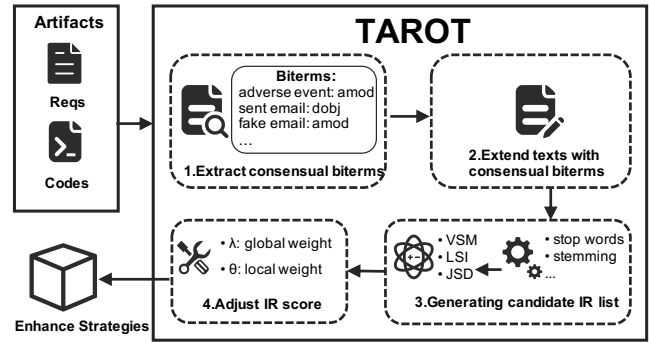


**Figure 2: Overview of the TAROT framework.**

recovery requires user verifications on the generated candidate links, a different body of work [18, 28, 52] uses user feedback (candidate links verified as either relevant links or false positives) to adjust the calculation of IR values. Panichella et al. [53] further combined the user feedback with code dependency analysis, and Gao et al. [23] proposed CLUSTER' that uses the closeness analysis on code dependencies [35] to improve IR-based approaches by propagating only a small amount of user feedback. Although the discussed approaches can improve the performance of IR-based traceability recovery, their provided benefits are still highly relevant to the quantity and quality of artifact texts, or the initially calculated IR values. Researchers tried to improve the quality of artifact texts by reducing the "noise" (discussed in Section 1), but the remaining texts are not enough to further improve IR-based approaches. In contrary, we use extracted consensual biterms to first enrich the corpus for IR techniques, and then adjust the ranking of candidate lists without considering code dependencies or user feedback. We then combined TAROT with CLUSTER' in our evaluation, and the experiment results showed that TAROT can collaborate with other enhancing strategies for further improvement.

**Using biterms in lexical analyses:** Biterms are first proposed for document retrieval [62] to address the data sparsity problem. For the same reason, Cheng et al. [13] use biterms to better establish topic models over short texts. Although the use of biterms in software engineering (SE) research is not many, researchers seem to agree that biterms are useful in analyzing short SE texts. For example, Hadi et al. [27] proposed an adaptive online biterm topic model to analyze version sensitive short texts. Instead of directly using the biterm topic modeling, we use biterms as the vital role of our approach to capture the same system functionalities that are described in different artifacts at different concept levels.

## 3 PROPOSED APPROACH

In this section, we introduce the proposed four-step approach TAROT. Figure 2 illustrates the overview of TAROT. First, we extract consensual biterms from both requirements and code. Second, we enrich artifact texts with consensual biterms (Section 3.2). Third, we generate candidate trace lists (Section 3.3). Fourth, we further adjust IR values with global and local weight, i.e., $\lambda$ and $\theta$ (Section 3.4). The updated candidate trace lists can also be the input for other enhancing strategies, e.g., our baseline CLUSTER' [23].
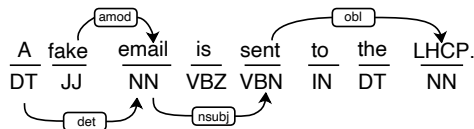
**Figure 3: Stanford typed dependencies parsing in a sentence**

## 3.1 Extracting Consensual Biterms

*3.1.1 Extracting Candidate Biterms from Requirements.* Two typical kinds of structured texts are often used to describe requirements in recent software development, i.e., use cases and repository issues. The use case usually consists of five parts, including title, precondition, main-flow, sub-flow, and alternative-flow, while the issue only has two parts: summary and description. We first split each requirement into several parts according to its text structure. Because requirements are written with natural language, we leverage Standford CoreNLP [42] to process each splitted part of text as follows: (1) we split text into independent sentences; (2) we tag each term's Part-of-speech (POS, e.g, nouns, verbs, adjectives, adverbs, and pronouns) in the sentence; (3) we conduct a Stanford typed dependencies parsing for each sentence that we can get the grammatical relationship between two terms in the sentence. For example, given a sentence as shown in Figure 3, each term's POS tag is labeled under itself. Stanford CoreNLP provides many types of tag, e.g., DT (determiner), JJ (adjective), NN (noun), VBZ (verb, third person singular present), VBN (verb, past tense), and IN (preposition or subordinating conjunction), etc. The grammatical relationship between two different terms is connected with arrowed lines, and relationships are labeled on the lines. For example, "fake" and "email" have an *amod* (adjectival modifier) relationship, "email" and "sent" have a *nsubj* (noun subject) relationship, and "sent" and "LHCP" (i.e., Licensed Health Care Professional) have a *obl* (oblique nominal) realtionship. Each two terms with grammatical relationship is recognized as a biterm (i.e., "**aEmail**", "**fakeEmail**", "**emailSent**", and "**sentLHCP**"). However, biterms are likely to carry different information values that depend on composed terms' POSs. On the one hand, previous studies [31, 41] have shown that verbs, nouns, and adjectives play important roles in expressing semantics. On the other hand, names of code entities (e.g., class name, method name, and variable name and type, etc.) are also generally composed of terms with the three POSs [49, 55]. Therefore, each term of a biterm should only be verb, noun, or adjective. Otherwise, the biterm will be ignored (e.g., "**aEmail**").

*3.1.2 Extracting Candidate Biterms from Code.* We extract biterms from both code identifier names (i.e., class name, method name, invoked method name, field name and its type, and parameter name and its type) and comments (i.e., comments for class or method). Unfortunately, existing POS taggers cannot achieved an accuracy of 100% in tagging code identifiers [3, 26, 49]. This is because code identifiers are not written in a natural language like requirements, so they do not contain the contexts of proper English sentences and grammatical structures. Nevertheless, to help developers better comprehend code, identifier names are usually composed of two

or more terms (usually no more than five terms [11, 49]) by following typical naming conventions, such as the `CamelCase` or the `snake_case`. Furthermore, identifier names tend to follow certain POS-based patterns among its splitted terms in general [10, 49, 55], such as **VB + JJ + NN** (e.g., getFakeEmail) and **VB + NN** (e.g., sendEmail) in `EmailUtil` shown in Figure 1. Therefore, we extract candidate biterms from identifier names by combining any two splitted terms sequentially. Take the identifier name getFakeEmail for example, we will extract three biterms including "**getFake**", "**getEmail**", and "**fakeEmail**". Although some unmeaning biterms may be extracted, i.e., "**getFake**", we consider them still tolerable because they are not many in one given identifier name. Furthermore, they are likely to be discarded when we further extract consensual biterms between requirements and code (discussed in Section 3.1.3). Meanwhile for code comments, we conduct the same candidate biterm extractions on them as those on requirements because they are both written in natural languages.

*3.1.3 Extracting Consensual Biterms.* We now use a quite straightforward crosscheck to extract the final consensual biterms from the two sets of candidate biterms from requirements and code, respectively. First, to address the unmeaning biterm problem from code candidate biterms (mentioned in Section 3.1.2), we only retain biterms that also appear in the set of requirement candidate biterms. The reason why we favor requirement candidate biterms is that we extract these biterms based on both the grammatical dependencies and the selected types of POSs (discussed in Section 3.1.1), so that the quality of requirement candidate biterms is likely to be much higher, let alone the typically better quality and quantity of requirement texts to guarantee the quality of Stanford CoreNLP output. Then we filtered out biterms that only appear in the set of requirement candidate biterms, to get the final set of consensual biterms between requirement and code. It is worth-while noting that each term in a given candidate biterm is restored to its original format, and the sequence of two terms is also omitted, except the domain abbreviations such as "LHCP" (which will be treated as a special noun). For example in our iTrust Sample, the "**emailSent**" mentioned in Section 3.1.1 and the "**sendEmail**" mentioned in Section 3.1.2 will both be saved as **(send, email)**, and this biterm becomes a consensual biterm after the crosscheck.

## 3.2 Enriching Texts with Consensual Biterms

*3.2.1 Enriching Requirements with Consensual Biterms.* To validly enrich requirement texts with extracted consensual biterms, we need to not only just "add them in the text", but also assign reasonable occurrence numbers for each added biterm so that our approach can effectively index these "artificial terms" later in calculating IR values between requirements and code classes. Specifically, we first record numbers of the occurrences of each consensual biterm respectively in different parts of a requirement. Then we accumulate the numbers in each part as the total number that the consensual biterms should be added into the requirement.

*3.2.2 Enriching Code with Consensual Biterms.* Similarly but differently, we design different count strategies according to the kind of identifier name in which the consensual biterm occurred. First, if a biterm appears in class names or method names, its count

pluses two. This is because the names of classes and methods are particularly important. Høst et al. [29] considered that "Methods are the smallest named units of aggregated behavior in most conventional programming languages and hence the cornerstone of abstraction". In addition, a class name is a high-level abstraction and summarization of a class, and a method name is a direct description of functionality. Then for a biterm appears in comments (i.e., comments of class and methods), its count pluses one for each occurrence. For invoked method name, field type and name, and parameter type and name, we take a conservative strategy. If the biterm only appears in these parts, its count will be assigned to one and not change no matter how many times it appears. Otherwise, we skip this biterm. Finally, we accumulate all counts of a biterm as the total number that should be added to each code class.

## 3.3 Generating Candidate Trace Links

**Creating Corpus**. Each class in code is extracted into one document containing its comments and identifiers including class names, method names, and field names. Then we use the CamelCase convention to split identifiers into terms. For each requirement, we extract a document that includes its title and content (e.g., preconditions, main-flow, sub-flows, and alternative-flows for structured use cases or title and description for commit issues).

**Normalizing Corpus**. We normalize the documents of requirements and classes by standard pre-processing techniques for IR including: (1) removing special tokens (e.g., punctuations); (2) converting all upper case letters into the lower case; (3) removing stop words; (4) performing the Porter stemming algorithm [56].

**Indexing Corpus and Computing Textual Similarities**. We use tf-idf for corpus indexing (for both terms and biterms) and three mainstream IR models to compute textual similarity: Vector Space Model (VSM) [6], Latent Semantic Indexing (LSI) [43], and the probabilistic Jensen and Shannon (JS) model [1].

**Generating Candidate Links**. We rank IR candidate lists in descending order based on the IR values of candidate links.

## 3.4 Adjusting IR Values through the Biterms

In addition to enriching artifact texts with biterms, we further leverage consensual biterms' global (i.e, for both a requirement and a code file) and local (i.e., for done part of a requirement) weights to adjust IR values to improve the ranking of generated IR candidate lists. Given a requirement $req$ and a code $cls$, we denote the set of biterms in $req$ and $cls$ as $BT_{req}$ and $BT_{cls}$ respectively. Then we can get the set of biterms that they shared denoted as

$$BT_{cons} = BT_{req} \cap BT_{cls}, \quad req \in UC \cup ISSUE, cls \in SC \quad (1)$$

where $UC$ and $ISSUE$ represent two different types of requirements, i.e., use cases and commit issues, respectively. $SC$ represents all source code. Note that we adopt each biterm's *inverse document frequency* (IDF) as its weight due to it can provide a good measure of biterms' uniqueness. For example, if a biterm occurred in almost all artifacts that are considered as a common biterm, then its IDF will be assigned a small value.

*3.4.1 Global Weight $\lambda$.* For a given requirement $req$ and code $cls$, we consider more biterms they shared, $req$ and $cls$ are more likely to be an actual trace link. In addition to the size of $BT_{cons}$, we also take

into account the proportion of $BT_{cons}$. We calculate proportions of $BT_{cons}$ in both $BT_{req}$ and $BT_{cls}$ respectively. Then, we combine these two proportions together and take half of the value as the global weight of $BT_{cons}$ denoted as $\lambda$:

$$\lambda(req, cls) = \left( \frac{\sum_{k=1}^{|BT_{cons}|} IDF_{bt_k}}{\sum_{i=1}^{|BT_{req}|} IDF_{bt_i}} + \frac{\sum_{k=1}^{|BT_{cons}|} IDF_{bt_k}}{\sum_{j=1}^{|BT_{cls}|} IDF_{bt_j}} \right) \times \frac{1}{2} \quad (2)$$

where $|BT_{req}|$, $|BT_{cls}|$, and $|BT_{cons}|$ represent the number of biterms in $req$, $cls$, and their shared biterms, respectively. $IDF_{bt_i}$, $IDF_{bt_j}$, and $IDF_{bt_k}$ are the IDFs of the $i_{th}$, $j_{th}$, and $k_{th}$ biterm in $BT_{req}$, $BT_{cls}$, and $BT_{cons}$, respectively.

*3.4.2 Local Weight $\theta$.* We further consider that each part of a requirement can differently contribute to the described system functionalities. We mainly consider two kinds of requirement artifacts: use cases and issues in software repositories. For use case, title and main-flow are summarized descriptions of the requirement. Therefore, they hold more information values and biterms in these two parts should be assigned higher weights, especially for title since it is more condensed. Sub-flows describe the main scenario involving detailed functionalities step-by-step. Alternative-flows describe alternate ways that the system behaves according to specific inputs. Precondition explains the state that the system must be in for the use case to be able to start which is not involved in functional description, so we do not consider this part. We then refer to previous study [4] and assign the following weights 0.4, 0.3, 0.2, and 0.1 to title, main-flow, sub-flow, and alternative-flow, respectively. Similarly, we assign weights 0.6 and 0.4 to summary and description, respectively for repository issues. Given a requirement $req$ and code $cls$, we compute the proportion of biterms they shared in all biterms that each part owns as follow:

$$\omega(req\_part, cls) = \frac{\sum_{k=1}^{|BT_{cons\_part}|} IDF_{bt_k}}{\sum_{i=1}^{|BT_{req\_part}|} IDF_{bt_i}} \quad (3)$$

where $req\_part$ is one part of $req$, $|BT_{cons\_part}|$ is the number of shared biterms for $req_{part}$ and $cls$. $IDF_{bt_k}$, and $IDF_{bt_i}$ are the IDFs of the $k_{th}$, $i_{th}$ in $BT_{req_part}$ and $BT_{cls}$ respectively. We then compute *local weight* (denoted as $\theta$) by combining weights from different part of requirement texts as follows:

$$\theta(req, cls) = \begin{cases} 0.4 \times \omega(title, cls) + 0.3 \times \omega(mf, cls) + \\ 0.2 \times \omega(sf, cls) + 0.1 \times \omega(af, cls), req \in UC \\ \\ 0.6 \times \omega(summ) + 0.4 \times \omega(desc), \quad req \in Issue \end{cases} \quad (4)$$

where $title$ refers to "title", $mf$ refers to "main-flow", $sf$ refers to "sub-flow", and $af$ refers to "alternative-flow" in use cases, while $summ$ is "summary" and $desc$ is "description" in issues.

Finally, we can leverage $\lambda$ and $\theta$ to adjust IR values if a requirement and code shared consensual biterms, i.e., $BT_{cons} \neq \emptyset$. Otherwise, we will penalize the IR values through multiplying origin IR value by 0.9 as follows:

$$IR_{new} = \begin{cases} IR_{initial} \times (1 + \lambda + \theta), & BT_{cons} \neq \emptyset \\ IR_{initial} \times 0.9, & BT_{cons} = \emptyset \end{cases} \quad (5)$$

where $IR_{initial}$ is the initial IR value and $IR_{new}$ is the updated one.

Table 1: Overview of the Nine Evaluated Systems

| # | iTrust | Gantt Project | Maven | Pig | Infinispan | Seam2 | Drools | Derby | Groovy |
|---|--------|---------------|-------|-----|------------|-------|--------|-------|--------|
| **Version** | 13.0 | 2.0.9 | 3.5.2 | 0.17.0 | 9.2.0 | 7.5.0 | 14.1.0 | 2.3.0 | 2.5.0 |
| **Programming Language** | Java | Java | Java | Java | Java | Java | Java | Java | Java |
| **KLoC** | 43 | 45 | 101 | 365 | 521 | 449 | 1091 | 136 | 381 |
| **Executed classes** | 137 | 124 | 82 | 289 | 319 | 150 | 248 | 611 | 100 |
| **Evaluated requirements** | 34 | 16 | 36 | 87 | 232 | 189 | 183 | 390 | 104 |
| **Relevant Traces in RTM** | 255 | 315 | 151 | 547 | 1116 | 463 | 841 | 2315 | 180 |

## 4 EXPERIMENT SETUP

We now introduce our experimental setup to evaluate our approach. Section 4.1 introduces nine evaluated systems. Section 4.2 defines metrics for evaluating the performance of our approach and baseline approaches. Section 4.3 discusses our research questions and the design of experiments.

### 4.1 Evaluated Systems

Our evaluation is based on nine real-world software systems: iTrust, GanttProject[5], Maven[6], Pig[7], Infinispan[8], Drools[9], Derby[10], Seam2[11], and Groovy[12]. Table 1 presents a summary of nine evaluated systems. We chose these systems because of their availability of both requirements specifications with developer maintained use cases and their Requirements-to-code Trace Matrices (RTM). For GanttProject, high-quality class-level traces are gained by recruiting the original developers. The RTM of iTrust contains method-level traces maintained by original developers and is publicly available. However, the RTMs of the other eight systems are class-level. To keep our experiment consistent at the same granularity, we propagated the method-level traces of iTrust to class-level traces by aggregating all traces to methods of a class on the class-level.

Meanwhile, Maven, Pig, Infinispan, Drools, Derby, Seam2, and Groovy come from the dataset named IlmSeven [59]. The dataset consists of seven open source projects where the traces between issues and changed classes for each project are elicited from its GitHub and the issue-tracking tool Jira. Specifically, Maven is the mainstream software project management tool. Pig is a platform to analyze large datasets consisting of a high-level language for expressing data analysis programs along with the infrastructure for assessing these programs. Infinispan is a distributed in-memory key-value NoSQL data store software. Drools is a business rule management system and reasoning engine for business policy and rules development, access, and change management. Derby is a relational database management system that can be embedded in Java programs and used for online transaction processing. Seam2 is a lightweight framework for building web applications in Java. Finally, Groovy is an object-oriented programming language for the Java platform. All projects are implemented in the Java language (the Groovy project using both Java and Groovy itself). We used all

seven projects and follow the same suggestions and heuristics proposed by Hui et al. [23] (also our baseline) to filter and merge each project's issues to avoid the likely noises and too fine-grained system functionalities when using the original issues as requirements. Our dataset is available at: https://github.com/huiAlex/TAROT.

### 4.2 Evaluation Metrics

We first leveraged two well-known metrics (i.e., recall and precision) to evaluate the performance of our approach. Precision represents the proportion of correct links among retrieved trace links, and recall represents the proportion of retrieved correct links among all correct links. They are defined as follows:

$$Precision = \frac{|relevant \cap retrieved|}{|retrieved|} \quad (6)$$

$$Recall = \frac{|relevant \cap retrieved|}{|relevant|} \quad (7)$$

where $relevant$ is the set of relevant links and $retrieved$ is the set of links retrieved by traceability recovery approaches.

A common way to evaluate the accuracy of IR techniques is to compare the precision values obtained at different recall levels, resulting in a set of precision-recall points displayed as curves. We then leveraged the following two metrics: AP (i.e., average precision) and MAP (i.e., mean average precision). AP and MAP are computed as:

$$AP = \frac{\sum_{r=1}^{N}(Precision(r) \times isRelevant(r))}{|RelevantDocuments|} \quad (8)$$

$$MAP = \frac{\sum_{q=1}^{Q} AP(q)}{Q} \quad (9)$$

where $r$ is the rank of the target artifact in an ordered list of links, $Precision(r)$ represents its precision value, $isRelevant()$ is a binary function assigned 1 if the link is relevant or 0 otherwise, $N$ is the number of all documents, $q$ is a single query, and $Q$ is the number of all queries. AP measures how well relevant documents of all queries (requirements) are ranked to the top of the retrieved links. Meanwhile, MAP uses the average of the AP scores of all queries to measure how well relevant documents for each query are ranked to the top of the retrieved links.

### 4.3 Research Questions

We aim to study whether the use of consensual biterms from text structures of requirements and code, i.e., enriching the input corpus and adjusting the calculated IR values, can help improve IR-based traceability recovery. Furthermore, because TAROT is proposed based on text structures only, we also want to study whether it can

collaborate with other enhancing strategies built upon different perspectives (e.g., code dependency analysis and the use of user feedback). Finally, we want to study how the corpus-enriching part and the value-adjusting part contribute separately to TAROT. Thus, we proposed the following three research questions:

**RQ1:** *Can TAROT help improve the performance of IR-based traceability recovery?*

**RQ2:** *Can TAROT collaborate with other enhancing strategies built upon different perspectives for IR-based traceability recovery?*

**RQ3:** *How many contributions do different parts of TAROT make individually?*

To study RQ1, we combine TAROT with the pure IR-based approach and named the combination as **IR-ONLY_TAROT**. We then compared the performance of IR-ONLY_TAROT and two baseline approaches. The first is the pure IR-based approach (**IR-ONLY**, without any enhancing strategies involved). Another one is **ConPOS**, which is the state-of-the-art approach that uses major POS categories and applies constraints to prune candidate IR links as a filtering process [2], while TAROT also uses POS taggers, but it chooses to first enrich the input corpora for IR techniques and then adjust the calculated IR values. We argue that TAROT have more potentials to better improving IR-based approaches because the quantity of either requirement texts or code texts is already quite small in practice. Correctly enriching those texts will be eventually more effective than further pruning them, and ConPOS is a really good baseline to evaluate our idea and approach. We then involved three mainstream IR models in our evaluation, i.e., VSM, LSI, and JS (which are widely used and evaluated in traceability research [15]), to compare IR-ONLY_TAROT with the two baseline approaches. Note that for the k value of LSI, we adopt a brute-force search strategy[33] to determine the optimal value of k for each evaluated system. For iTrust, GanttProject, Maven, Pig, Infinispan, Seam2, Drools, Derby, and Groovy, their calibrated k values are 85, 65, 70, 140, 170, 180, 300, 130, and 175, respectively. For new systems, we suggest using the same calibration process before fine-tuning them to optimize the performance.

In addition to the metrics proposed in Section 4.2, we also use a statistical significance test to check whether the performance of IR-ONLY_TAROT is significantly better than that of baseline approaches. Through the significance test used in reference [5, 34, 35], we use the F-measure at each recall point as the single dependent variable of our significant test. We use the F-measure because we want to know whether IR-ONLY_TAROT improves both accuracy and recall. F-measured value is calculated as follows:

$$F = \frac{2P \times R}{P + R} \tag{10}$$

where $P$ represents precision, $R$ represents recall, and $F$ is the harmonic mean of $P$ and $R$. A higher F-measure means that both precision and recall are high and balanced. We then use the Wilcoxon rank sum test [64] to test the following null hypothesis:

$H_0$: *There is no difference between the performance of TAROT and baseline approaches.*

We use $\alpha = 0.05$ to accept or refute the null hypothesis $H_0$. We also use a non-parametric effect size measure for ordinal data, i.e.,

Cliff's *delta* ($\delta$) [38], to quantify the amount of difference between IR-ONLY_TAROT and two baseline appproaches as follows:

$$\delta = \left| \frac{\# (x_1 > x_2) - \# (x_1 < x_2)}{n_1 n_2} \right| \tag{11}$$

where $x_1$ and $x_2$ represent F-measures of TAROT and a chosen baseline approach, and $n_1$ and $n_2$ are the sizes of the sample groups. The effect size is considered negligible for $\delta < 0.15$, small for $0.15 \leq \delta < 0.33$, medium for $0.33 \leq \delta < 0.47$, large for $\delta > 0.47$.

To study RQ2, we choose **CLUSTER'** as the representative enhancing strategies built upon different perspectives for IR-based traceability recovery to test whether TAROT can collaborate with other enhancements. CLUSTER' propagates a small amount of user feedback with closeness analysis on code dependencies [23]. We use IR-ONLY_TAROT to generate IR candidate links and make them as the input of CLUSTER'. We name this combination of TAROT and CLUSTER' as **CLUSTER'_TAROT**. Like RQ1, we also compared CLUSTER' and CLUSTER'_TAROT on VSM, LSI, and JS, use F-measure at each recall point to see whether CLUSTER'_TAROT can improve both accuracy and recall, use the Wilcoxon rank sum test to test $H_0$ ($\alpha = 0.05$), and use Cliff's *delta* to quantify the amount of difference between CLUSTER' and CLUSTER'_TAROT.

To study RQ3, we perform an ablation study to evaluate the performance of TAROT's different parts. We regard IR-ONLY as the basics and add each part of TAROT to the basics incrementally. By comparing the performance of the approach before and after adding a part can see how many contributions that each step of TAROT makes. Specifically, "+*b*" represents when we only use extracted consensual biterms to improve the input quality for IR techniques, without considering the *global weight* $\lambda$ and the *local weight* $\theta$. "+ $\lambda$" represents when we adjust IR values through the biterms, the *global weight* $\lambda$ will be taken into consideration. "+ $\theta$" represents when we adjust IR values through the biterms, the *local weight* $\theta$ will be taken into consideration. We also use VSM, LSI, and JS, to compare IR-ONLYs with different steps included.

## 5 DISCUSSIONS ON RESULTS

**RQ1: Can TAROT help improve the performance of IR-based traceability recovery?** Table 2 shows the experiment results of nine evaluated systems (rows). For each system and each IR model (columns), we compared the performance of IR-ONLY, ConPOS, and IR-ONLY_TAROT. Sub column 1 shows the average precision (AP). Sub column 2 shows the mean average precision (MAP). Sub column 3 shows the *p*-value of the F-measure significance test for IR-ONLY_TAROT and sub column 4 shows the Cliff's $\delta$. For each approach, the best result of AP, MAP, and *p*-value $\leq 0.05$ are in bold text. The results show that IR-ONLY_TAROT outperforms ConPOS on both AP and MAP for all cases. When compared with IR-ONLY, IR-ONLY_TAROT performs better on MAP for all case, only performs worse on AP in 3 cases of DroolsAP on Drools (0.12 on average). Specifically, in 42 out of 54 cases, the F-measure for the result of IR-ONLY_TAROT is significantly higher than those of the two baseline approaches (p-value $\leq 0.05$) at each level of recall, which indicates that IR-ONLY_TAROT outperforms baseline approaches in the majority of cases. Figure 4 shows and compares

**Table 2: The number of computed AP, MAP, p-value, and Cliff's $\delta$ evaluating each approach (evaluated systems iTrust, GanttProject, Maven, Pig, Infnispan, Seam2, Drools, Derby, and Groovy combined with IR models VSM, LSI, and JS)**

| | | VSM | | | | LSI | | | | JS | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AP | MAP | $p$-value | C'$\delta$ | AP | MAP | $p$-value | C'$\delta$ | AP | MAP | $p$-value | C'$\delta$ |
| iTrust | IR-ONLY | 45.78 | 58.43 | **0.05** | 0.10 | 46.01 | 59.17 | 0.06 | 0.10 | 40.57 | 56.01 | **<0.01** | 0.14 |
| | ConPOS | 46.89 | 59.03 | 0.13 | 0.08 | 46.60 | 59.00 | 0.17 | 0.07 | 39.95 | 55.27 | **0.02** | 0.12 |
| | IR-ONLY_TAROT | **49.50** | **62.12** | - | - | **49.18** | **61.50** | - | - | **45.74** | **58.59** | - | - |
| | CLUSTER' | 64.08 | 68.51 | 0.54 | 0.03 | 64.55 | 69.94 | 0.34 | 0.05 | 61.75 | 66.53 | **0.03** | 0.11 |
| | CLUSTER'_TAROT | **65.90** | **72.33** | - | - | **67.02** | **72.95** | - | - | **64.20** | **69.09** | - | - |
| GanttProject | IR-ONLY | 43.17 | 49.79 | **<0.01** | 0.14 | 43.89 | 51.72 | **<0.01** | 0.15 | 36.50 | 46.76 | **<0.01** | 0.22 |
| | ConPOS | 44.33 | 51.41 | **<0.01** | 0.17 | 44.81 | 53.15 | **<0.01** | 0.18 | 38.24 | 50.04 | **<0.01** | 0.24 |
| | IR-ONLY_TAROT | **47.37** | **54.09** | - | - | **48.63** | **55.25** | - | - | **43.90** | **51.76** | - | - |
| | CLUSTER' | 62.49 | 71.90 | **<0.01** | 0.18 | 66.28 | **72.13** | 0.09 | 0.08 | 70.69 | 74.10 | 0.97 | 0.0 |
| | CLUSTER'_TAROT | **73.18** | **75.38** | - | - | **70.91** | 72.05 | - | - | **73.19** | **75.37** | - | - |
| Maven | IR-ONLY | 22.27 | 37.11 | **<0.01** | 0.24 | 22.15 | 40.65 | **<0.01** | 0.20 | 24.08 | 40.29 | **<0.01** | 0.29 |
| | ConPOS | 22.34 | 37.12 | **<0.01** | 0.17 | 21.23 | 38.67 | **<0.01** | 0.18 | 22.56 | 37.33 | **<0.01** | 0.24 |
| | IR-ONLY_TAROT | **26.77** | **43.65** | - | - | **26.97** | **47.17** | - | - | **25.50** | **48.35** | - | - |
| | CLUSTER' | 42.53 | 47.31 | 0.61 | 0.03 | 41.92 | 49.61 | 0.75 | 0.02 | 42.00 | 49.95 | 0.21 | 0.08 |
| | CLUSTER'_TAROT | **45.17** | **53.91** | - | - | **44.18** | **54.87** | - | - | **47.55** | **55.75** | - | - |
| Pig | IR-ONLY | 19.71 | 36.37 | **0.02** | 0.08 | 17.89 | 36.62 | 0.12 | 0.05 | 14.64 | 31.90 | **<0.01** | 0.13 |
| | ConPOS | 20.05 | 39.18 | **0.03** | 0.07 | 17.96 | 37.67 | 0.13 | 0.05 | 13.46 | 32.94 | **<0.01** | 0.13 |
| | IR-ONLY_TAROT | **22.93** | **39.38** | - | - | **21.09** | **37.88** | - | - | **20.51** | **37.58** | - | - |
| | CLUSTER' | **31.11** | **43.56** | **<0.01** | 0.12 | **29.49** | **40.44** | 0.16 | 0.05 | 27.39 | 36.84 | **<0.01** | 0.10 |
| | CLUSTER'_TAROT | 26.47 | 42.98 | - | - | 28.31 | 39.60 | - | - | **28.96** | **41.26** | - | - |
| Infinispan | IR-ONLY | 8.73 | 25.44 | **<0.01** | 0.09 | 9.05 | 26.84 | **<0.01** | 0.10 | 7.32 | 26.02 | **<0.01** | 0.14 |
| | ConPOS | 9.45 | 26.70 | **<0.01** | 0.08 | 9.48 | 26.86 | **<0.01** | 0.09 | 7.68 | 26.61 | **<0.01** | 0.13 |
| | IR-ONLY_TAROT | **11.34** | **29.04** | - | - | **12.00** | **30.68** | - | - | **10.17** | **27.82** | - | - |
| | CLUSTER' | 18.77 | 30.50 | 0.13 | 0.04 | **20.70** | 32.64 | 0.35 | 0.02 | 20.92 | **31.82** | **0.02** | 0.06 |
| | CLUSTER'_TAROT | **19.73** | **32.61** | - | - | 20.69 | **34.37** | - | - | **21.34** | 31.12 | - | - |
| Seam2 | IR-ONLY | 18.99 | 40.61 | **<0.01** | 0.12 | 17.19 | 42.49 | **<0.01** | 0.10 | 16.64 | 40.47 | **<0.01** | 0.11 |
| | ConPOS | 20.26 | 41.96 | **<0.01** | 0.11 | 18.65 | 43.49 | **<0.01** | 0.10 | 17.33 | 42.01 | **<0.01** | 0.11 |
| | IR-ONLY_TAROT | **23.65** | **46.11** | - | - | **21.85** | **46.41** | - | - | **20.84** | **43.06** | - | - |
| | CLUSTER' | 34.16 | 49.32 | 0.08 | 0.07 | 33.45 | 47.47 | **0.02** | 0.09 | 30.32 | 45.95 | **0.02** | 0.09 |
| | CLUSTER'_TAROT | **36.14** | **52.71** | - | - | **37.06** | **52.54** | - | - | **34.6** | **51.43** | - | - |
| Drools | IR-ONLY | **8.87** | 21.06 | 0.12 | 0.04 | **7.98** | 20.98 | **0.02** | 0.07 | **7.56** | 21.20 | 0.14 | 0.04 |
| | ConPOS | 8.48 | 21.27 | 0.13 | 0.04 | 7.03 | 21.41 | **0.02** | 0.07 | 6.67 | 21.25 | 0.10 | 0.05 |
| | IR-ONLY_TAROT | 8.83 | **22.41** | - | - | 7.84 | **22.46** | - | - | 7.39 | **21.86** | - | - |
| | CLUSTER' | **15.86** | 24.65 | **0.02** | 0.07 | 13.62 | 24.46 | **<0.01** | 0.07 | **15.08** | 23.88 | 0.63 | 0.01 |
| | CLUSTER'_TAROT | 15.77 | **26.72** | - | - | **14.60** | **26.24** | - | - | 13.08 | **24.24** | - | - |
| Derby | IR-ONLY | 10.23 | 26.95 | **<0.01** | 0.20 | 9.58 | 26.74 | **<0.01** | 0.17 | 9.32 | 26.21 | **<0.01** | 0.18 |
| | ConPOS | 11.58 | 30.32 | **<0.01** | 0.13 | 10.31 | 29.76 | **<0.01** | 0.12 | 8.38 | 27.69 | **<0.01** | 0.20 |
| | IR-ONLY_TAROT | **15.05** | **35.07** | - | - | **13.29** | **35.29** | - | - | **14.22** | **34.01** | - | - |
| | CLUSTER' | 24.81 | 36.66 | **<0.01** | 0.06 | 24.11 | 34.35 | **<0.01** | 0.16 | 23.67 | 32.71 | **<0.01** | 0.13 |
| | CLUSTER'_TAROT | **25.05** | **43.51** | - | - | **27.91** | **42.23** | - | - | **26.23** | **41.42** | - | - |
| Groovy | IR-ONLY | 26.98 | 52.72 | **0.02** | 0.15 | 27.40 | 55.50 | **0.02** | 0.14 | 19.69 | 44.02 | **<0.01** | 0.22 |
| | ConPOS | 28.19 | 56.39 | 0.47 | 0.04 | 27.83 | 58.99 | 0.52 | 0.04 | 17.23 | 48.61 | 0.89 | 0.01 |
| | IR-ONLY_TAROT | **33.59** | **58.96** | - | - | **33.53** | **61.80** | - | - | **27.91** | **54.44** | - | - |
| | CLUSTER' | 41.12 | 69.75 | **<0.01** | 0.27 | 43.93 | 73.69 | **<0.01** | 0.30 | 36.86 | 64.82 | **<0.01** | 0.34 |
| | CLUSTER'_TAROT | **46.49** | **70.83** | - | - | **48.52** | **73.74** | - | - | **46.70** | **69.53** | - | - |

the precision-recall curves for the three approaches grouped by each system and IR model.

We now use the adapted excerpt from the iTrust system to demonstrate why IR-ONLY_TAROT is able to outperform IR-ONLY and ConPOS. For use case UC35 and code `EmailUtil` as shown in Figure 1, they shared two biterms, i.e., "**sendEmail**" and "**fakeEmail**". Because the two biterms both appears in sub-flows of UC35 twice, TAROT adds them into UC35 twice respectively according to the biterms extension rules (discussed in Section 3.2.1). For `EmailUtil`, TAROT adds "**sendEmail**" three times (appears in class comment and method name) and "**fakeEmail**" once (appears in invoked method name). It turns out that TAROT can successfully enrich

texts of UC35 and `EmailUtil`. Furthermore, TAROT further uses *global weight* $\lambda$ and *local weight* $\theta$ to improve the IR value of "UC35 -> EmailUtil". For the result of IR-ONLY, the ranks of `EmailUtil` in the candidate list of UC35 and lists of all 34 use cases are 40/137 and 1172/4658 respectively, while IR-ONLY_TAROT promotes them to 8/137 and 201/4658 respectively. Yet for ConPOS, there is only one verb shared by UC35 and `EmailUtil`, i.e., "send", which does not help to achieve much improvement.

**RQ2: Can TAROT collaborate with other enhancing strategies for IR-based traceability recovery?** Table 2 shows the performances of CLUSTER' and CLUSTER' _TAROT on nine evaluated
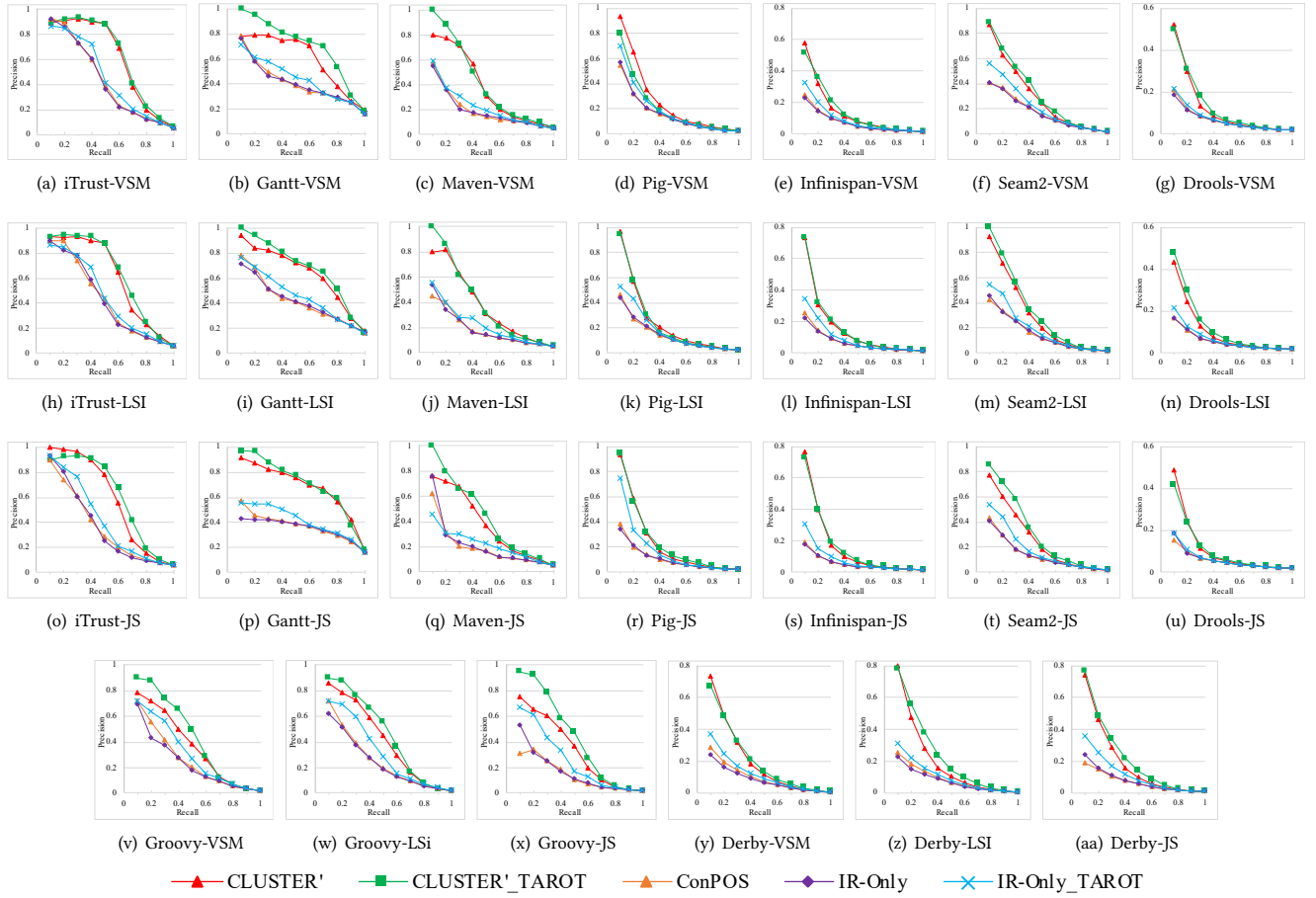
**Figure 4: Precision/Recall curves grouped by evaluated systems (subgraph (a), (h), (o) for iTrust, and similarly for Gantt, Maven, Pig, Infinispan, Seam2 and Drools, while (v), (w), (x) for Groovy, and similarly for Derby) and IR models (VSM, LSI, and JS).**

systems. In 15 out of 27 cases, the F-measure for the result of CLUSTER' _TAROT is significantly higher than that of CLUSTER' (p-value ≤ 0.05) at each level of recall, which indicates that CLUSTER' _TAROT outperforms CLUSTER' in the majority of cases. When compared with CLUSTER', CLUSTER' _TAROT outperforms in 22 out of 27 cases on AP. Its AP is 2.49 higher than that of CLUSTER' on average in total 27 cases. The highest value is 10.69 on GanttProject-VSM. Meanwhile, CLUSTER' _TAROT outperforms CLUSTER' in 23 out of 27 cases on MAP. Its MAP is 3.16 higher than that of CLUSTER' on average in total 27 cases. The highest value is 8.71 on Derby-JS. It turns out that TAROT can collaborate with CLUSTER', which is a totally different IR enhancing strategies built upon both code dependency analysis and the use of user feedback. The combined CLUSTER'_TAROT can not only achieve further improvement where CLUSTER' has already achieves good results (e.g., iTrust and GanttProject), but also achieve indeed improvement where CLUSTER' achieves relatively low results (e.g., Infinispan and Drools). Finding more effective way to collaborate TAROT with other enhancing strategies is one of our future work. This result

also inspires us to further explore whether TAROT can be beneficial to ML-based techniques in our future work, because TAROT inherently improves the quality and quantity of requirement and code texts.

**RQ3: How many contributions does different parts of TAROT make individually?** Table 3 shows the experiment results of nine evaluated systems (rows). For each system and each IR model (columns), we compared the performance of incrementally adding each part of TAROT including "+ $b$" means only extending artifact texts with consensual biterms on the bases of the pure IR-based approach (IR-ONLY), "+ $b$ + $\lambda$" means further introducing global weight $\lambda$, and "+ $b$ + $\lambda$ + $\theta$" means the entire TAROT. Sub column 1 shows the average precision (AP). Sub column 2 shows the mean average precision (MAP). For each approach, the best result of AP, MAP, and $p$-value ≤ 0.05 are in bold text. From the table, we can observe that "+$b$" outperforms IR-ONLY in all 27 cases on AP (on average 1.69) and in 24 out of 27 cases on MAP (on average 2.18), which indicates that TAROT can extract funcationality consensual biterms and promote the quality of artifact texts. In addition, the results show that "+ $b$ + $\lambda$" outperforms "+$b$" in 24 out of 27 cases on

**Table 3: The number of computed AP and MAP evaluating each part of TAROT ($b$ = consensual biterms, $\lambda$ = gloabal weight, and $\theta$ = local weight)**

| | | VSM | | LSI | | JS | |
|---|---|---|---|---|---|---|---|
| | | AP | MAP | AP | MAP | AP | MAP |
| iTrust | IR-ONLY | 45.78 | 58.43 | 46.01 | 59.17 | 40.57 | 56.01 |
| | + $b$ | 48.28 | 58.43 | 47.85 | 59.26 | 44.66 | 58.40 |
| | + $b$ + $\lambda$ | 49.16 | 60.49 | 49.04 | 59.72 | 45.89 | 58.71 |
| | + $b$ + $\lambda$ + $\theta$ | 49.50 | 62.12 | 49.18 | 61.50 | 45.74 | 58.59 |
| GanttProject | IR-ONLY | 43.17 | 49.79 | 43.89 | 51.72 | 36.50 | 46.76 |
| | + $b$ | 45.72 | 53.28 | 46.19 | 53.24 | 41.23 | 50.81 |
| | + $b$ + $\lambda$ | 46.19 | 53.26 | 46.92 | 54.43 | 41.38 | 50.76 |
| | + $b$ + $\lambda$ + $\theta$ | 47.37 | 54.09 | 48.63 | 55.25 | 43.90 | 51.76 |
| Maven | IR-ONLY | 22.27 | 37.11 | 22.15 | 40.65 | 24.08 | 40.29 |
| | + $b$ | 24.42 | 41.47 | 22.75 | 45.27 | 25.81 | 43.39 |
| | + $b$ + $\lambda$ | 26.64 | 44.22 | 25.60 | 46.81 | 26.01 | 47.05 |
| | + $b$ + $\lambda$ + $\theta$ | 26.77 | 43.65 | 26.97 | 47.17 | 25.50 | 48.35 |
| Pig | IR-ONLY | 19.71 | 36.37 | 17.89 | 36.62 | 14.64 | 31.90 |
| | + $b$ | 21.55 | 37.75 | 19.52 | 36.92 | 17.70 | 34.35 |
| | + $b$ + $\lambda$ | 22.23 | 37.61 | 20.41 | 36.60 | 18.44 | 34.98 |
| | + $b$ + $\lambda$ + $\theta$ | 22.93 | 39.38 | 21.09 | 37.88 | 20.51 | 37.58 |
| Infinispan | IR-ONLY | 8.73 | 25.44 | 9.05 | 26.84 | 7.32 | 26.02 |
| | + $b$ | 9.76 | 27.53 | 9.99 | 28.72 | 8.25 | 27.81 |
| | + $b$ + $\lambda$ | 10.71 | 28.29 | 11.05 | 29.75 | 9.40 | 27.66 |
| | + $b$ + $\lambda$ + $\theta$ | 11.34 | 29.04 | 12.00 | 30.68 | 10.17 | 27.82 |
| Seam2 | IR-ONLY | 18.99 | 40.61 | 17.19 | 42.49 | 16.64 | 40.47 |
| | + $b$ | 21.58 | 43.25 | 17.81 | 42.70 | 17.71 | 39.95 |
| | + $b$ + $\lambda$ | 24.05 | 44.76 | 21.34 | 45.14 | 20.53 | 42.26 |
| | + $b$ + $\lambda$ + $\theta$ | 23.65 | 46.11 | 21.85 | 46.41 | 20.84 | 43.06 |
| Drools | IR-ONLY | 8.87 | 21.06 | 7.98 | 20.98 | 7.56 | 21.20 |
| | + $b$ | 9.12 | 22.57 | 8.75 | 23.10 | 7.68 | 21.43 |
| | + $b$ + $\lambda$ | 8.99 | 22.10 | 8.45 | 23.37 | 7.47 | 21.54 |
| | + $b$ + $\lambda$ + $\theta$ | 8.83 | 22.41 | 7.84 | 22.46 | 7.39 | 21.86 |
| Derby | IR-ONLY | 10.23 | 26.95 | 9.58 | 26.74 | 9.32 | 26.21 |
| | + $b$ | 12.42 | 30.83 | 10.18 | 28.88 | 10.63 | 28.50 |
| | + $b$ + $\lambda$ | 13.78 | 33.21 | 11.77 | 32.64 | 12.26 | 31.71 |
| | + $b$ + $\lambda$ + $\theta$ | 15.05 | 35.07 | 13.29 | 35.29 | 14.22 | 34.01 |
| Groovy | IR-ONLY | 26.98 | 52.72 | 27.40 | 55.50 | 19.69 | 44.02 |
| | + $b$ | 27.94 | 52.27 | 28.08 | 54.69 | 22.15 | 46.52 |
| | + $b$ + $\lambda$ | 31.34 | 55.99 | 31.75 | 58.11 | 26.93 | 52.94 |
| | + $b$ + $\lambda$ + $\theta$ | 33.59 | 58.96 | 33.53 | 61.80 | 27.91 | 54.44 |

AP (on average 1.69) and 21 out of 27 cases on MAP (on average 2.02) respectively, which means shared consensual biterms can narrow down the semantic gap. Moreover, the results show that "$b + \lambda + \theta$" can outperform "+ $b + \lambda$" in 21 out of 27 cases on AP (on average 1.13) and 24 out of 27 cases on MAP (on average 1.43) respectively. We also noted that all three parts of TAROT can outperform IR-ONLY in almost all cases on both AP and MAP, only in 4 out of 162 cases performs worse.

The overall evaluation results also implied that TAROT is likely to be beneficial in practice because: (1) it can improve IR-based traceability by solely analyzing the text structures without requiring additional analysis (such as code instrumentation, which may not be viable for real-world projects); (2) it can still work even on systems with really low-quality texts (e.g., Infinispan, Drools, and Derby in our evaluation); (3) it can complement other enhancing strategies, and is also likely to be useful for ML-based traceability recovery because TAROT inherently improves the quality and quantity of requirement and code texts.

## 6    THREATS TO VALIDITY

**Internal threats.** One possible internal threat to our approach is that we cannot guarantee 100% accuracy in segmenting texts, tagging POSs, and parsing dependencies based on Stanford CoreNLP. However, existing work has reported that the accuracy of off-the-shelf NLP tools is acceptable when analyzing texts with the context

of proper sentences and grammatical structures [3], and we found no obvious errors in the output of Stanford CoreNLP as well. We will consider SE-specific NLP tools (e.g., [54]) in future work.

**External threats.** Our evaluation uses the same dataset of CLUSTER' [23] with nine systems including seven open-source systems. We think that our findings from this evaluation are relevant because these evaluated systems are either widely studied or used in practice from different domains [34, 59]. Furthermore, we combined the evaluated systems with three mainstream IR models (i.e., VSM, LSI, and JS) to extend our experiments to a total of 27 variants (e.g., iTrust-VSM and Pig-JS). However, one possible threat is that our evaluation is only based on part of the system functionalities of Maven, Pig, Infinispan, Drools, Derby, Seam2, and Groovy because Gao et al. set up heuristics to filter and merge issues with linked commit logs from the IlmSeven dataset to elicit both requirements and RTMs, thus covering only part of the extracted issues from Jira and GitHub. Nevertheless, we use the same RTMs to compare our approach with the baseline approaches on each evaluated system, thus making no bias.

## 7    CONCLUSIONS AND FUTURE WORK

In this study, we propose to extract co-occurred word pairs from the text structures of both requirements and code (i.e., consensual biterms) to improve IR-based traceability recovery. The consensual biterms are extracted by our crosscheck between candidate requirement biterms and candidate code biterms. We argue that although the consensual biterms are just correlative term combinations that can be located in both requirements and code, they represents important semantics of the system that are consistently followed during the software development process, thus being valuable for automated traceability recovery. We then use these biterms to first enrich the input corpora for IR techniques, and then adjust the generated IR values. An empirical evaluation based on nine real-world systems shows that our approach can not only outperform baseline approaches, but also collaborate with other enhancing strategies built upon different perspectives. In future work we first plan to extract more consensual biterms from additional software artifacts (e.g., design and tests). We then plan to study whether consensual biterms can improve ML-based traceability recovery as well, such as improving IR values used as features in ML-based approaches, or enhancing word-embeddings with additional biterms. Our replication package is available at: https://github.com/huiAlex/TAROT.

Hui Gao, Hongyu Kuang, Kexin Sun, Xiaoxing Ma, Alexander Egyed, Patrick Mäder, Guoping Rong, Dong Shao, and He Zhang

# REFERENCES

[1] Ahron Abadi, Mordechai Nisenson, and Yahalomit Simionovici. 2008. A Traceability Technique for Specifications. In *The 16th IEEE International Conference on Program Comprehension, ICPC 2008, Amsterdam, The Netherlands, June 10-13, 2008*, René L. Krikhaar, Ralf Lämmel, and Chris Verhoef (Eds.). IEEE Computer Society, 103–112.

[2] Nasir Ali, Haipeng Cai, Abdelwahab Hamou-Lhadj, and Jameleddine Hassine. 2018. Exploiting Parts-of-Speech for Effective Automated Requirements Traceability. *Information and Software Technology* 106 (09 2018). https://doi.org/10.1016/j.infsof.2018.09.009

[3] Nasir Ali, Haipeng Cai, Abdelwahab Hamou-Lhadj, and Jameleddine Hassine. 2019. Exploiting Parts-of-Speech for effective automated requirements traceability. *Inf. Softw. Technol.* 106 (2019), 126–141. https://doi.org/10.1016/j.infsof.2018.09.009

[4] Nasir Ali, Zohreh Sharafi, Yann-Gaël Guéhéneuc, and Giuliano Antoniol. 2012. An empirical study on requirements traceability using eye-tracking. In *28th IEEE International Conference on Software Maintenance, ICSM 2012, Trento, Italy, September 23-28, 2012*. IEEE Computer Society, 191–200. https://doi.org/10.1109/ICSM.2012.6405271

[5] Nasir Ali, Zohreh Sharafi, Yann-Gaël Guéhéneuc, and Giuliano Antoniol. 2015. An empirical study on the importance of source code entities for requirements traceability. *Empirical Software Engineering* 20, 2 (2015), 442–478.

[6] Giuliano Antoniol, Gerardo Canfora, Gerardo Casazza, Andrea De Lucia, and Ettore Merlo. 2002. Recovering Traceability Links between Code and Documentation. *IEEE Trans. Software Eng.* 28, 10 (2002), 970–983.

[7] Ricardo Baezayates and Berthier Ribeironeto. 2011. *Modern information retrieval*. Addison-Wesley Publishing CompanyUnited States.

[8] Ted J. Biggerstaff, Bharat G. Mitbander, and Dallas E. Webster. 1993. The Concept Assignment Problem in Program Understanding. In *Proceedings of the 15th International Conference on Software Engineering, Baltimore, Maryland, USA, May 17-21, 1993*, Victor R. Basili, Richard A. DeMillo, and Takuya Katayama (Eds.). IEEE Computer Society / ACM Press, 482–498. http://portal.acm.org/citation.cfm?id=257572.257679

[9] D. Binkley. 2007. Source Code Analysis: A Road Map. In *Future of Software Engineering (FOSE '07)*. IEEE Computer Society, Los Alamitos, CA, USA, 104–119.

[10] David W. Binkley, Matthew Hearn, and Dawn J. Lawrie. 2011. Improving identifier informativeness using part of speech information. In *Proceedings of the 8th International Working Conference on Mining Software Repositories, MSR 2011 (Co-located with ICSE), Waikiki, Honolulu, HI, USA, May 21-28, 2011, Proceedings*, Arie van Deursen, Tao Xie, and Thomas Zimmermann (Eds.). ACM, 203–206. https://doi.org/10.1145/1985441.1985471

[11] Simon Butler, Michel Wermelinger, and Yijun Yu. 2015. A survey of the forms of Java reference names. In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension, ICPC 2015, Florence/Firenze, Italy, May 16-24, 2015*, Andrea De Lucia, Christian Bird, and Rocco Oliveto (Eds.). IEEE Computer Society, 196–206. https://doi.org/10.1109/ICPC.2015.30

[12] Giovanni Capobianco, Andrea De Lucia, Rocco Oliveto, Annibale Panichella, and Sebastiano Panichella. 2013. Improving IR-based traceability recovery via noun-based indexing of software artifacts. *J. Softw. Evol. Process.* 25, 7 (2013), 743–762. https://doi.org/10.1002/smr.1564

[13] Xueqi Cheng, Xiaohui Yan, Yanyan Lan, and Jiafeng Guo. 2014. BTM: Topic Modeling over Short Texts. *IEEE Transactions on Knowledge and Data Engineering* 26, 12 (2014), 2928–2941. https://doi.org/10.1109/TKDE.2014.2313872

[14] Jane Cleland-Huang. 2013. Are Requirements Alive and Kicking? *IEEE Softw.* 30, 3 (2013), 13–15.

[15] Jane Cleland-Huang, Orlena Gotel, Jane Huffman Hayes, Patrick Mäder, and Andrea Zisman. 2014. Software traceability: trends and future directions. In *Proceedings of the on Future of Software Engineering, FOSE 2014, Hyderabad, India, May 31 - June 7, 2014*, James D. Herbsleb and Matthew B. Dwyer (Eds.). ACM, 55–69.

[16] Jane Cleland-Huang, Raffaella Settimi, Chuan Duan, and Xuchang Zou. 2005. Utilizing Supporting Evidence to Improve Dynamic Requirements Traceability. In *13th IEEE International Conference on Requirements Engineering (RE 2005), 29 August - 2 September 2005, Paris, France*. IEEE Computer Society, 135–144.

[17] Andrea De Lucia, Andrian Marcus, Rocco Oliveto, and Denys Poshyvanyk. 2012. Information Retrieval Methods for Automated Traceability Recovery. In *Software and Systems Traceability*, Jane Cleland-Huang, Olly Gotel, and Andrea Zisman (Eds.). Springer, 71–98.

[18] Andrea De Lucia, Rocco Oliveto, and Paola Sgueglia. 2006. Incremental Approach and User Feedbacks: a Silver Bullet for Traceability Recovery. In *22nd IEEE International Conference on Software Maintenance (ICSM 2006), 24-27 September 2006, Philadelphia, Pennsylvania, USA*. IEEE Computer Society, 299–309.

[19] Andrea De Lucia, Massimiliano Di Penta, Rocco Oliveto, Annibale Panichella, and Sebastiano Panichella. 2011. Improving IR-based Traceability Recovery Using Smoothing Filters. In *The 19th IEEE International Conference on Program Comprehension, ICPC 2011, Kingston, ON, Canada, June 22-24, 2011*. IEEE Computer Society, 21–30.

[20] Diana Diaz, Gabriele Bavota, Andrian Marcus, Rocco Oliveto, Silvia Takahashi, and Andrea De Lucia. 2013. Using code ownership to improve IR-based Traceability Link Recovery. In *IEEE 21st International Conference on Program Comprehension, ICPC 2013, San Francisco, CA, USA, 20-21 May, 2013*. IEEE Computer Society, 123–132. https://doi.org/10.1109/ICPC.2013.6613840

[21] Bogdan Dit, Meghan Revelle, and Denys Poshyvanyk. 2013. Integrating information retrieval, execution and link analysis algorithms to improve feature location in software. *Empirical Software Engineering* 18, 2 (2013), 277–309.

[22] Alexander Egyed, Florian Graf, and Paul Grünbacher. 2010. Effort and Quality of Recovering Requirements-to-Code Traces: Two Exploratory Experiments. In *RE 2010, 18th IEEE International Requirements Engineering Conference, Sydney, New South Wales, Australia, September 27 - October 1, 2010*. IEEE Computer Society, 221–230.

[23] Hui Gao, Hongyu Kuang, Xiaoxing Ma, Hao Hu, Jian Lü, Patrick Mäder, and Alexander Egyed. 2022. Propagating frugal user feedback through closeness of code dependencies to improve IR-based traceability recovery. *Empir. Softw. Eng.* 27, 2 (2022), 41. https://doi.org/10.1007/s10664-021-10091-5

[24] Malcom Gethers, Rocco Oliveto, Denys Poshyvanyk, and Andrea De Lucia. 2011. On integrating orthogonal information retrieval methods to improve traceability recovery. In *IEEE 27th International Conference on Software Maintenance, ICSM 2011, Williamsburg, VA, USA, September 25-30, 2011*. IEEE Computer Society, 133–142.

[25] Jin Guo, Jinghui Cheng, and Jane Cleland-Huang. 2017. Semantically enhanced software traceability using deep learning techniques. In *Proceedings of the 39th International Conference on Software Engineering, ICSE 2017, Buenos Aires, Argentina, May 20-28, 2017*, Sebastián Uchitel, Alessandro Orso, and Martin P. Robillard (Eds.). IEEE / ACM, 3–14.

[26] Samir Gupta, Sana Malik, Lori L. Pollock, and K. Vijay-Shanker. 2013. Part-of-speech tagging of program identifiers for improved text-based software engineering tools. In *IEEE 21st International Conference on Program Comprehension, ICPC 2013, San Francisco, CA, USA, 20-21 May, 2013*. IEEE Computer Society, 3–12. https://doi.org/10.1109/ICPC.2013.6613828

[27] Mohammad Abdul Hadi and Fatemeh H Fard. 2020. AOBTM: Adaptive Online Biterm Topic Modeling for Version Sensitive Short-texts Analysis. In *2020 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. 593–604. https://doi.org/10.1109/ICSME46990.2020.00062

[28] Jane Huffman Hayes, Alex Dekhtyar, and Senthil Karthikeyan Sundaram. 2006. Advancing Candidate Link Generation for Requirements Tracing: The Study of Methods. *IEEE Trans. Software Eng.* 32, 1 (2006), 4–19.

[29] Einar W. Høst and Bjarte M. Østvold. 2009. Debugging Method Names. In *ECOOP 2009 - Object-Oriented Programming, 23rd European Conference, Genoa, Italy, July 6-10, 2009. Proceedings (Lecture Notes in Computer Science, Vol. 5653)*, Sophia Drossopoulou (Ed.). Springer, 294–317. https://doi.org/10.1007/978-3-642-03013-0_14

[30] Daniel Jurafsky and James Martin. 2008. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Vol. 2.

[31] Reza Karimpour, Amineh Ghorbani, Azadeh Pishdad, Mitra Mohtarami, Abolfazl AleAhmad, Hadi Amiri, and Farhad Oroumchian. 2008. Improving Persian Information Retrieval Systems Using Stemming and Part of Speech Tagging. In *Evaluating Systems for Multilingual and Multimodal Information Access, 9th Workshop of the Cross-Language Evaluation Forum, CLEF 2008, Aarhus, Denmark, September 17-19, 2008, Revised Selected Papers (Lecture Notes in Computer Science, Vol. 5706)*, Carol Peters, Thomas Deselaers, Nicola Ferro, Julio Gonzalo, Gareth J. F. Jones, Mikko Kurimo, Thomas Mandl, Anselmo Peñas, and Vivien Petras (Eds.). Springer, 89–96. https://doi.org/10.1007/978-3-642-04447-2_10

[32] Edward Keenan. 1975. *Formal Semantics of Natural Language*. https://doi.org/10.1017/CBO9780511897696

[33] Saket Khatiwada, Miroslav Tushev, and Anas Mahmoud. 2020. On Combining IR Methods to Improve Bug Localization. In *ICPC '20: 28th International Conference on Program Comprehension, Seoul, Republic of Korea, July 13-15, 2020*. ACM, 252–262. https://doi.org/10.1145/3387904.3389280

[34] Hongyu Kuang, Hui Gao, Hao Hu, Xiaoxing Ma, Jian Lu, Patrick Mäder, and Alexander Egyed. 2019. Using frugal user feedback with closeness analysis on code to improve IR-based traceability recovery. In *Proceedings of the 27th International Conference on Program Comprehension, ICPC 2019, Montreal, QC, Canada, May 25-31, 2019*, Yann-Gaël Guéhéneuc, Foutse Khomh, and Federica Sarro (Eds.). IEEE / ACM, 369–379.

[35] Hongyu Kuang, Jia Nie, Hao Hu, Patrick Rempel, Jian Lu, Alexander Egyed, and Patrick Mäder. 2017. Analyzing closeness of code dependencies for improving IR-based Traceability Recovery. In *IEEE 24th International Conference on Software Analysis, Evolution and Reengineering, SANER 2017, Klagenfurt, Austria, February 20-24, 2017*, Martin Pinzger, Gabriele Bavota, and Andrian Marcus (Eds.). IEEE Computer Society, 68–78.

[36] Tien-Duy B. Le, Mario Linares Vásquez, David Lo, and Denys Poshyvanyk. 2015. RCLinker: automated linking of issue reports and commits leveraging rich contextual information. In *Proceedings of the 2015 IEEE 23rd International Conference*

on Program Comprehension, ICPC 2015, Florence/Firenze, Italy, May 16-24, 2015, Andrea De Lucia, Christian Bird, and Rocco Oliveto (Eds.). IEEE Computer Society, 36–47. https://doi.org/10.1109/ICPC.2015.13

[37] Jinfeng Lin, Yalin Liu, Qingkai Zeng, Meng Jiang, and Jane Cleland-Huang. 2021. Traceability Transformed: Generating more Accurate Links with Pre-Trained BERT Models. In 43rd IEEE/ACM International Conference on Software Engineering, ICSE 2021, Madrid, Spain, 22-30 May 2021. IEEE, 324–335. https://doi.org/10.1109/ICSE43902.2021.00040

[38] Guillermo Macbeth, Eugenia Razumiejczyk, and Rubén Daniel Ledesma. 2011. Cliff's Delta Calculator: A non-parametric effect size program for two groups of observations. Universitas Psychologica 10, 2 (2011), 545–555.

[39] Patrick Mäder and Alexander Egyed. 2012. Assessing the effect of requirements traceability for software maintenance. In 28th IEEE International Conference on Software Maintenance, ICSM 2012, Trento, Italy, September 23-28, 2012. IEEE Computer Society, 171–180.

[40] Patrick Mäder, Paul L. Jones, Yi Zhang, and Jane Cleland-Huang. 2013. Strategic Traceability for Safety-Critical Projects. IEEE Softw. 30, 3 (2013), 58–66. https://doi.org/10.1109/MS.2013.60

[41] Anas Mahmoud and Nan Niu. 2015. On the role of semantics in automated requirements tracing. Requir. Eng. 20, 3 (2015), 281–300. https://doi.org/10.1007/s00766-013-0199-y

[42] Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, System Demonstrations. The Association for Computer Linguistics, 55–60. https://doi.org/10.3115/v1/p14-5010

[43] Andrian Marcus and Jonathan I. Maletic. 2003. Recovering Documentation-to-Source-Code Traceability Links using Latent Semantic Indexing. In Proceedings of the 25th International Conference on Software Engineering, May 3-10, 2003, Portland, Oregon, USA, Lori A. Clarke, Laurie Dillon, and Walter F. Tichy (Eds.). IEEE Computer Society, 125–137.

[44] Collin McMillan, Denys Poshyvanyk, and Meghan Revelle. 2009. Combining textual and structural analysis of software artifacts for traceability link recovery. In ICSE Workshop on Traceability in Emerging Forms of Software Engineering, TEFSE@ICSE 2009. Vancouver, BC, Canada, 18 May, 2009, Giuliano Antoniol, Denys Poshyvanyk, and Rocco Oliveto (Eds.). IEEE Computer Society, 41–48.

[45] Chris Mills, Javier Escobar-Avila, Aditya Bhattacharya, Grigoriy Kondyukov, Shayok Chakraborty, and Sonia Haiduc. 2019. Tracing with Less Data: Active Learning for Classification-Based Traceability Link Recovery. In 2019 IEEE International Conference on Software Maintenance and Evolution, ICSME 2019, Cleveland, OH, USA, September 29 - October 4, 2019. IEEE, 103–113. https://doi.org/10.1109/ICSME.2019.00020

[46] Chris Mills, Javier Escobar-Avila, and Sonia Haiduc. 2018. Automatic Traceability Maintenance via Machine Learning Classification. In 2018 IEEE International Conference on Software Maintenance and Evolution, ICSME 2018, Madrid, Spain, September 23-29, 2018. IEEE Computer Society, 369–380. https://doi.org/10.1109/ICSME.2018.00045

[47] Kevin Moran, David N. Palacio, Carlos Bernal-Cárdenas, Daniel McCrystal, Denys Poshyvanyk, Chris Shenefiel, and Jeff Johnson. 2020. Improving the effectiveness of traceability link recovery using hierarchical bayesian networks. In ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020, Gregg Rothermel and Doo-Hwan Bae (Eds.). ACM, 873–885. https://doi.org/10.1145/3377811.3380418

[48] Shiva Nejati, Mehrdad Sabetzadeh, Davide Falessi, Lionel C. Briand, and Thierry Coq. 2012. A SysML-based approach to traceability management and design slicing in support of safety certification: Framework, tool support, and case studies. Inf. Softw. Technol. 54, 6 (2012), 569–590. https://doi.org/10.1016/j.infsof.2012.01.005

[49] Christian D. Newman, Reem S. Alsuhaibani, Michael John Decker, Anthony Peruma, Dishant Kaushik, Mohamed Wiem Mkaouer, and Emily Hill. 2020. On the generation, structure, and semantics of grammar patterns in source code identifiers. J. Syst. Softw. 170 (2020), 110740. https://doi.org/10.1016/j.jss.2020.110740

[50] Armstrong Nhlabatsi, Yijun Yu, Andrea Zisman, Thein Than Tun, Niamul Khan, Arosha K. Bandara, Khaled M. Khan, and Bashar Nuseibeh. 2015. Managing Security Control Assumptions Using Causal Traceability. In 8th IEEE/ACM International Symposium on Software and Systems Traceability, SST 2015, Florence, Italy, May 17, 2015, Patrick Mäder and Rocco Oliveto (Eds.). IEEE Computer Society, 43–49. https://doi.org/10.1109/SST.2015.14

[51] Annibale Panichella, Bogdan Dit, Rocco Oliveto, Massimilano Di Penta, Denys Poshynanyk, and Andrea De Lucia. 2013. How to effectively use topic models for software engineering tasks? An approach based on Genetic Algorithms. In 2013 35th International Conference on Software Engineering (ICSE). 522–531. https://doi.org/10.1109/ICSE.2013.6606598

[52] Annibale Panichella, Andrea De Lucia, and Andy Zaidman. 2015. Adaptive User Feedback for IR-Based Traceability Recovery. In 8th IEEE/ACM International Symposium on Software and Systems Traceability, SST 2015, Florence, Italy, May 17, 2015, Patrick Mäder and Rocco Oliveto (Eds.). IEEE Computer Society, 15–21.

[53] Annibale Panichella, Collin McMillan, Evan Moritz, Davide Palmieri, Rocco Oliveto, Denys Poshyvanyk, and Andrea De Lucia. 2013. When and How Using Structural Information to Improve IR-Based Traceability Recovery. In 17th European Conference on Software Maintenance and Reengineering, CSMR 2013, Genova, Italy, March 5-8, 2013, Anthony Cleve, Filippo Ricca, and Maura Cerioli (Eds.). IEEE Computer Society, 199–208.

[54] Profir-Petru Pârtachi, Santanu Kumar Dash, Christoph Treude, and Earl T. Barr. 2020. POSIT: simultaneously tagging natural and programming languages. In ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020, Gregg Rothermel and Doo-Hwan Bae (Eds.). ACM, 1348–1358. https://doi.org/10.1145/3377811.3380440

[55] Anthony Peruma, Emily Hu, Jiajun Chen, Eman Abdullah AlOmar, Mohamed Wiem Mkaouer, and Christian D. Newman. 2021. Using Grammar Patterns to Interpret Test Method Name Evolution. In 29th IEEE/ACM International Conference on Program Comprehension, ICPC 2021, Madrid, Spain, May 20-21, 2021. IEEE, 335–346. https://doi.org/10.1109/ICPC52881.2021.00039

[56] Martin F. Porter. 1980. An algorithm for suffix stripping. Program 14, 3 (1980), 130–137. https://doi.org/10.1108/eb046814

[57] Denys Poshyvanyk, Yann-Gaël Guéhéneuc, Andrian Marcus, Giuliano Antoniol, and Václav Rajlich. 2007. Feature Location Using Probabilistic Ranking of Methods Based on Execution Scenarios and Information Retrieval. IEEE Trans. Software Eng. 33, 6 (2007), 420–432.

[58] Balasubramaniam Ramesh and Matthias Jarke. 2001. Toward Reference Models of Requirements Traceability. IEEE Trans. Software Eng. 27, 1 (2001), 58–93. https://doi.org/10.1109/32.895989

[59] Michael Rath, Patrick Rempel, and Patrick Mäder. 2017. The IlmSeven Dataset. In 25th IEEE International Requirements Engineering Conference, RE 2017, Lisbon, Portugal, September 4-8, 2017, Ana Moreira, João Araújo, Jane Hayes, and Barbara Paech (Eds.). IEEE Computer Society, 516–519.

[60] Michael Rath, Jacob Rendall, Jin L. C. Guo, Jane Cleland-Huang, and Patrick Mäder. 2019. Traceability in the Wild: Automatically Augmenting Incomplete Trace Links. In Software Engineering and Software Management, SE/SWM 2019, Stuttgart, Germany, February 18-22, 2019 (LNI, Vol. P-292), Steffen Becker, Ivan Bogicevic, Georg Herzwurm, and Stefan Wagner (Eds.). GI, 63.

[61] Patrick Rempel and Patrick Mäder. 2017. Preventing Defects: The Impact of Requirements Traceability Completeness on Software Quality. IEEE Trans. Software Eng. 43, 8 (2017), 777–797.

[62] Munirathnam Srikanth and Rohini Srihari. 2002. Biterm Language Models for Document Retrieval. In Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (Tampere, Finland) (SIGIR '02). Association for Computing Machinery, New York, NY, USA, 425–426. https://doi.org/10.1145/564376.564476

[63] Yan Sun, Qing Wang, and Ye Yang. 2017. FRLink: Improving the recovery of missing issue-commit links by revisiting file relevance. Inf. Softw. Technol. 84 (2017), 33–47. https://doi.org/10.1016/j.infsof.2016.11.010

[64] Frank Wilcoxon. 1944. Individual Comparisons by Ranking Methods. Biom Bull. Biometrics 1, 6 (1944), 80–83.