

# Automatic Traceability Maintenance via Machine Learning Classification

Chris Mills, Javier Escobar-Avila, Sonia Haiduc  
Department of Computer Science - Florida State University  
Tallahassee, FL, USA  
{cmills, escobara, shaiduc}@cs.fsu.edu

**Abstract**—Previous studies have shown that *software traceability*, the ability to link together related artifacts from different sources within a project (e.g., source code, use cases, documentation, etc.), improves project outcomes by assisting developers and other stakeholders with common tasks such as impact analysis, concept location, etc. Establishing traceability links in a software system is an important and costly task, but only half the struggle. As the project undergoes maintenance and evolution, new artifacts are added and existing ones are changed, resulting in outdated traceability information. Therefore, specific steps need to be taken to make sure that traceability links are maintained in tandem with the rest of the project.

In this paper we address this problem and propose a novel approach called TRAIL for maintaining traceability information in a system. The novelty of TRAIL stands in the fact that it leverages previously captured knowledge about project traceability to train a machine learning classifier which can then be used to derive new traceability links and update existing ones. We evaluated TRAIL on 11 commonly used traceability datasets from six software systems and compared it to seven popular Information Retrieval (IR) techniques including the most common approaches used in previous work. The results indicate that TRAIL outperforms all IR approaches in terms of precision, recall, and F-score.

## I. INTRODUCTION

Software systems are comprised of information stored in various artifacts such as source code, bug reports, requirements specifications, use and test cases, interaction diagrams, and user documentation among others. Traceability Link Recovery (TLR) is the software engineering task focused on establishing bidirectional links (i.e., *traceability links*) between related artifacts of different types in order to provide developers and stakeholders a detailed picture of how a system is constructed and the relationships between system components. The resulting software traceability naturally supports other tasks such as concept location, impact analysis, program comprehension, verifying test coverage, ensuring that system and regulatory requirements are met, etc. and has been proven to be useful in practice [1], [2], [3].

Establishing traceability links between the artifacts of a system is extremely arduous and error-prone when performed manually. This has led to a large body of research proposing techniques that aid developers with this task [4], [5], [6], [7], [8]. However, establishing traceability links at one point in the lifetime of a software system is only part of the struggle. As the system evolves over time, some artifacts get deleted, new artifacts are added that have not yet been linked to others, and substantial changes can break existing links. Therefore, the

benefits of software traceability can quickly degrade unless this information is updated and maintained in tandem with the evolving software artifacts [9], [10], [11].

In this paper we address this problem and propose a novel approach for traceability maintenance. Our approach, called TRAIL (TRAcability lInk cLassifier), uses historically collected traceability information (i.e., existing traceability links between pairs of artifacts) to train a machine learning classifier which is then able to classify the link between any new or existing pair of artifacts as *valid* (i.e., the two artifacts are related) or *invalid* (i.e., the two artifacts are unrelated). Since TRAIL relies only on the features of the two artifacts in order to determine the validity of the link between them, it is able to classify links between brand new artifacts introduced in the system, as well as to reassess existing links, when the artifacts involved have changed. To the best of our knowledge TRAIL is the first approach that requires neither human intervention nor a set of predefined rules for these tasks.

We evaluate TRAIL on a set of 11 datasets from six software systems which are commonly used in traceability studies [12], [13], [14], [15], [16]. We also compare TRAIL to Information Retrieval (IR) techniques, which are the most popular type of technique used in traceability link recovery and maintenance [4]. In particular, we compare TRAIL to seven IR approaches, including the most common approaches used in previous work [5], [6], [7], [17]. The results of our evaluation reveal that TRAIL is able to achieve high precision, recall and F-measure, significantly outperforming all IR approaches.

The main contributions of this paper are:

- 1) A novel approach to traceability maintenance called TRAIL which uses historical trace information to train a machine learning classification algorithm that predicts if a new or updated traceability link is valid or invalid.
- 2) An empirical derivation of TRAIL that shows which feature selection, balancing technique, and classification algorithm provide the best performance.
- 3) An empirical evaluation of TRAIL on 11 popular traceability datasets, which indicates that TRAIL is able to achieve high precision, recall, and F-measure.
- 4) A comparison of TRAIL with seven IR approaches previously used in traceability link recovery and maintenance, which shows that TRAIL significantly outperforms IR in terms of precision, recall, and F-measure.

The remainder of the paper is structured as follows: section II presents our approach; sections III and IV describe the evaluation we performed and a discussion of its results; section V discusses threats to validity; section VI provides an overview of related work; and section VIII summarizes our conclusions and presents ideas for future research.

## II. APPROACH

In this section we introduce our novel approach to traceability link maintenance, called TRAIL (TRAcability lInk cLassifier). Different from existing approaches, TRAIL leverages historical traceability information (i.e., pre-existing traceability links) to infer how artifacts should be linked. TRAIL transforms the traceability maintenance problem into a binary classification problem where the links between pairs of artifacts are classified as being valid or invalid. It employs machine learning algorithms that use features derived from the existing traceability links to infer statistical patterns that differentiate the valid and invalid links between artifacts. Due to its construction, TRAIL can assess both existing and new pairs artifacts in the system, which represents an advantage over previous approaches.

Rather than setting in stone the algorithms, features, and configurations, we designed TRAIL as a framework where individual components can be replaced by others that perform a similar task. This ensures that TRAIL is adaptable and configurable to specific needs that could arise in other domains or applications. In the following subsections we first present the framework itself at a high level, and then describe the specific implementation we considered in our study.

The TRAIL framework consists of four main components:

- 1) A set of *features* that are meant to represent the potential traceability links between software artifacts;
- 2) A *feature selection* component that extracts representative features from all the ones initially considered, in order to reduce the dimensionality, avoid overfitting, and improve the learning and generalization of the classifier.
- 3) A *data balancing* component that, if needed, rebalances the training data between the valid and invalid classes of links, to ensure better model training and performance;
- 4) A *machine learning classification algorithm* that is used to determine which potential links are valid and which are invalid.

Based on the above components, a few steps need to be performed in order to instantiate and use TRAIL. First, feature engineering is employed to establish a set of features that can be used to represent potential traceability links (section II-A describes our features). Then a feature selection algorithm is selected and applied to extract only the most relevant features (more details in section II-B). Next, because the approach represents all possible links between two sets of artifacts, the training data is expected to be highly imbalanced (i.e., given two sets of artifacts, there are usually much fewer *valid* links than *invalid* ones between them). To address this situation, a rebalancing technique is selected and applied to the data prior to training the classifier; this is further discussed in section II-C. Finally, a classification algorithm is chosen and trained with

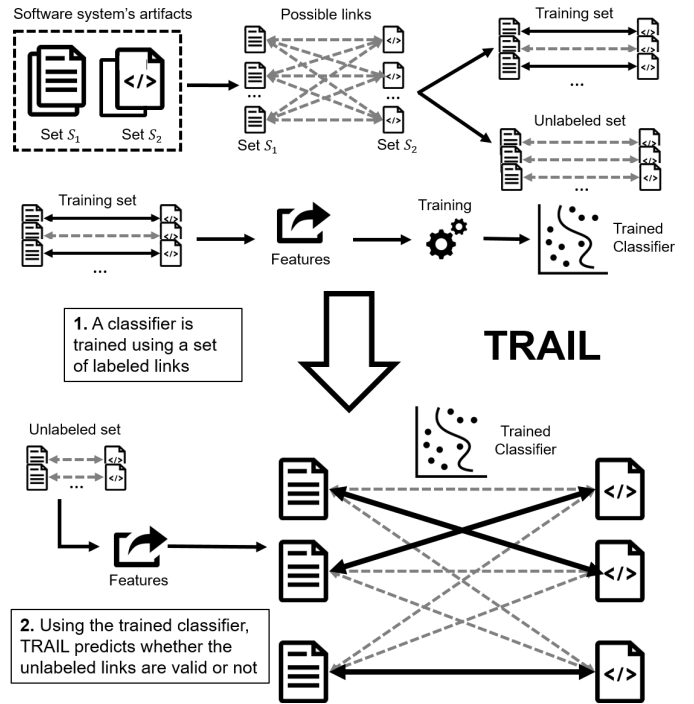


Fig. 1. The general TRAIL framework for traceability link classification the resulting data and can then be used to classify new links. Figure 1 shows an overview of how an instance of TRAIL is trained and used to predict the validity of traceability links.

Note that when implementing TRAIL the implementor is at liberty to design a variation that meets her specific needs by choosing a completely custom set of features, feature selection algorithm, rebalancing technique, and classifier. Moreover, TRAIL can easily be extended with emerging techniques in machine learning, which allows it to implicitly benefit from a rapidly evolving discipline and research area. In the subsections that follow, we introduce the specific settings we used to instantiate TRAIL for the study in this paper.

### A. Features Representing the Links

TRAIL considers all of the possible links that could exist between two given sets of software artifacts and predicts each to be either valid (i.e., the artifacts are related) or invalid (i.e., the artifacts are not related). Formally, if we have two artifact sets  $S_1$  and  $S_2$ , the approach predicts the validity of each element in the Cartesian product  $S_1 \times S_2$ . That is, for each artifact  $s \in S_1$ , we predict the validity of the potential link that exists between  $s$  and some  $s' \in S_2$ . In order to perform this classification, we must have a set of features that define an internal, vector representation for our classifier. The features used in our implementations of TRAIL can be separated into three categories: IR-based, query quality (QQ), and document statistics features.

1) *IR-based Features*: While IR-based techniques do not provide a silver bullet for traceability, their ability to capture the semantic similarity between software artifacts has long been established. Therefore, while our approach does not use

IR to propose lists of candidate links to a human, we still leverage the power of IR in the first set of features.

Given two artifact sets  $S_1$  and  $S_2$ , and a possible link between artifacts  $s1 \in S_1$  and  $s2 \in S_2$ , we capture the strength of this link based on IR using two metrics. First, we use  $s1$  as a query and the artifacts in  $S_2$  as the corpus. After running  $s1$  as a query through an IR engine, we capture the rank at which  $s2$  appears in the list of results as the first metric. Then we repeat the procedure, this time considering  $s2$  as the query,  $S_1$  as the corpus, and capturing the rank of  $s1$  in the list of results as the second metric. Therefore, the representation of each potential link in TRAIL considers trace link recovery from either direction (i.e., using either artifact in the link as a query), whereas traditional IR approaches consider only a single direction. This is an important distinction, as previous work has shown that the choice of retrieval direction impacts retrieval performance specifically for traceability [18].

In this study, we use seven different IR approaches to compute each of the two aforementioned metrics, resulting in a total of *14 different IR-based features for each possible link*. The IR approaches we use are: Vector Space Model with TF-IDF weighting and cosine similarity, Okapi BM25, Jensen Shannon, Latent Semantic Analysis, Latent Dirichlet Allocation, and two approaches based on smoothing methods for language models: Dirichlet and Jelenik-Mercer [19]. Next we provide a short description of each of these approaches.

*Vector Space Model (VSM)*: In the VSM, software artifacts/documents are represented using a *term-by-document* matrix. Each element of the matrix stores the importance of a term in the document and corpus, expressed using its *term frequency-inverse document frequency* (TF-IDF). Since each document is represented as a vector, we can compute the similarity between documents using cosine similarity:

$$sim_{cosine}(d_1, d_2) = \frac{\sum_{i=1}^t d_{1,i} d_{2,i}}{\sum_{i=1}^t d_{1,i}^2 \cdot \sum_{i=1}^t d_{2,i}^2} \quad (1)$$

where  $d_1$  and  $d_2$  are two documents in the corpus,  $t$  is the total number of unique terms in the corpus,  $d_{1,i}$  and  $d_{2,i}$  are the TF-IDF weights of the  $i$ th-term in each document.

*Latent Semantic Analysis (LSA)*: LSA [20] is able to capture information about the co-occurrence of terms, addressing the synonymy and polysemy problems that VSM is not able to. LSA uses Singular Value Decomposition to decompose the term-by-document matrix into three matrices, one of them containing a matrix of singular values ( $S_0$ ). By taking the  $k$  largest values from  $S_0$ , a new term-by-document matrix can be reconstructed with low dimensionality and information about the association between terms. The document representation in the reconstructed matrix together with cosine similarity is used to compute similarity between documents.

*Latent Dirichlet Allocation (LDA)*: LDA is a generative probabilistic model that represents each document as a mixture of latent topics, and each topic as a distribution over words in the corpus [21]. After using LDA, each document is represented as a vector of probabilities, each one describing the probability of a topic to appear in the document. We use this representation,

together with Hellinger distance, to compute the similarity between documents. Specifically, we define the similarity between two LDA document representations  $d_1$  and  $d_2$  as:

$$sim_{LDA}(d_1, d_2) = 1 - \frac{1}{\sqrt{2}} \|(\sqrt{d_1} - \sqrt{d_2})\| \quad (2)$$

Since our goal is to accurately estimate the similarity between two documents in the topic space, we configured the LDA implementation<sup>1</sup> with a large number of topics (250).

*Jensen-Shannon (JS)*: The JS model also represents software artifacts as probability distributions over terms in the corpus via hypothesis testing techniques [6]. Given the JS representation of two documents  $d_1$  and  $d_2$ , their similarity is computed as follows:

$$sim_{JS}(d_1, d_2) = 1 - \left[ H\left(\frac{d_1 + d_2}{2}\right) - \frac{H(d_1) + H(d_2)}{2} \right] \quad (3)$$

$$H(d) = \sum_{w \in W} -P(w) \cdot \log P(w)$$

where  $w$  is a word in document  $d$ ,  $W$  is the set of unique words in document  $d$ , and  $P(w)$  is the probability of word  $w$  appearing in document  $d$ .

*Okapi BM25*: The BM25 model scores each document in the corpus based on the query terms appearing in it. The scoring function is the following [22]:

$$score_{BM25}(q, d) = \frac{\left( \sum_{t \in q} \log \left[ \frac{N}{df_t} \right] \right) \cdot (\lambda + 1) \cdot tf_{t,d}}{tf_{t,d} \cdot \lambda \left( (1 - b) + b \cdot \left( \frac{L_d}{L_{ave}} \right) \right)} \quad (4)$$

where  $q$  is the query,  $d$  is a document in the corpus,  $N$  is the number of documents in the corpus,  $df_t$  is the number of documents the term  $t$  appears in,  $tf_{t,d}$  is the term frequency of term  $t$  in document  $d$ ,  $L_d$  is the length of document  $d$  (in number of words),  $L_{ave}$  is the average document length in the corpus,  $b$  is a parameter used to control how much effect field-length normalization should have, and  $\lambda$  is a positive parameter that calibrates the document term frequency scaling. We used Lucene's<sup>2</sup> default implementation of BM25, which utilizes  $\lambda = 1.2$  and  $b = 0.75$ .

*Language Model with Dirichlet (LM-Dirichlet) and Jelinek-Mercer smoothing (LM-JM)*: These two models first define a language model for each document in the corpus, and then rank the documents according to the probability that the language model of each document  $d$  generates a query  $q$ . Both models use *smoothing* to improve accuracy by adjusting the maximum estimator of a language model[19].

A general form of a smoothed language model is the following:

$$p(q, d) = \prod_{w \in q} p(q_i|d) \quad (5)$$

where

<sup>1</sup>We use the LDA implementation offered by the R package *topicmodels*.

<sup>2</sup><https://lucene.apache.org/> Version 6.4.1

$$p(q_i|d) = \begin{cases} p_{smooth}(q_i|d) & \text{if word } q_i \text{ is seen in } d \\ \alpha \cdot p(q_i|C) & \text{otherwise} \end{cases}$$

$\alpha$  is a coefficient that controls the probability assigned to unseen words, and  $C$  is the entire corpus. Due to space restrictions, we do not discuss the implementation details of LM-Dirichlet or LM-JM; however, they can be found in [19].

2) *Query Quality Features*: We also use a set of metrics that were applied in previous work as estimators for the quality of software artifacts when used as queries for IR [23], [24], and were subsequently applied specifically to the task of identifying hard-to-trace artifacts in TLR [25]. We use these metrics as they complement our IR-based features and can give the classifier more contextual information about the link between two artifacts. For example, if the IR rank of an artifact in a potential link is low (i.e., a poor textual match with the other artifact in the link), query quality (QQ) metrics can give an indication of whether this is due to the artifact being generally hard-to-trace, or to the fact that the two artifacts are indeed not related (i.e., linked).

We adopted all 28 QQ metrics used in previous work in software engineering [25]. The metrics can be split into two main categories: pre-retrieval (21 metrics), which can be applied without running the query and capture general properties of the text found in the artifact, and post-retrieval (7 metrics), which also take into account the ranked list of results returned when running the artifact as a query. Each main category further contains sub-categories, which focus on different properties of a document. Due to space constraints, we refer the interested reader to a previous study that used these metrics for concept location and TLR [25].

For each possible link, we applied each pre-retrieval metric to the two documents in the potential link (since each one can be used as a query), resulting in *42 different pre-retrieval QQ features for each link*. Further, we computed each of the seven post-retrieval metrics using five different IR approaches and two different retrieval directions (similarly to the IR-based features described above, post-retrieval QQ metrics also depend on the retrieval direction), resulting in *70 different post-retrieval QQ features for each link*. Note that we do not compute post-retrieval features for LSA or LDA as most of these metrics require some form of document perturbation, which results in the need to re-index the space in which documents are represented for each query. This is particularly cost intensive for these two approaches. Further, document perturbations are performed many times per metric to minimize the effect of non-determinism, which makes it computationally infeasible to use these metrics as features for our representation.

3) *Document Statistics Features*: In addition to the aforementioned attributes, we also include some simple document-level features. These metrics are intended to gauge document relevance through term overlap and provide information on the size of documents as a proxy for the information contained in the document. The three document features we use are: the number of unique terms in a document, the total number of

terms in a document (including duplicates), and the percentage of overlapping terms between the two documents in a candidate link. Therefore, we compute two features for each artifact in a possible link as well as one feature for the link itself, resulting in a total of *5 features for each link*.

In summary, we extracted 131 features from each potential traceability link: 14 IR-based features, 42 pre-retrieval QQ metrics, 70 post-retrieval QQ metrics, and 5 document features. We normalized each of these to the interval [0,1].

## B. Feature Selection

Each feature included in a predictive model’s internal representation contains some possibility for error. One mechanism to limit the impact of this potential error is to construct a so-called *parsimonious* model: one that explains a phenomena with the lowest dimensional internal representation possible. Parsimonious models also limit the effects of spurious correlations, reduce the probability of overfitting, and minimize computation time required to generate the representation. There are many techniques for lowering the dimensionality of a feature space. In this work we consider five algorithms implemented in Weka. The first is Correlation-based Feature Subset Selection (CFS-Subset), which considers the predictive power of each independent variable as well as their mutual correlation. The remaining four algorithms consider the relationship between each independent variable and the dependent (i.e., outcome) variable in terms of: Pearson’s Correlation, Gain Ratio, Information Gain, and Symmetrical Uncertainty.

## C. Class Imbalance and Data Rebalancing

In any given software project, it is expected that the number of possible traceability links is much larger than the number of valid links that actually exist between related artifacts. As such, the boolean classification performed by TRAIL has an inherent *class imbalance*, which can make it difficult to differentiate minority class instances due to insufficient data for learning a pattern [26]. Indeed, if there are only two valid links in a total set of 100 possible links, a classifier can optimize accuracy by predicting all links as invalid, unfortunately misclassifying the two links that are of the highest importance. In the context of our classification, this translates to low recall, which severely impacts the applicability of an automated approach to traceability maintenance.

To address this problem, rebalancing techniques can be applied to data used to train the classifier, which provides a more equivalent statistical representation of the majority and minority classes. Undersampling and oversampling are two widely used approaches to deal with class imbalance [26]. Undersampling involves reducing the majority class by selecting only a subset of its datapoints for training. Oversampling tries to increase the number of datapoints in the minority class, often by artificially creating new datapoints based on those in the original data. We employ both types of sampling methods in our evaluation in order to determine if they help improve the results of TRAIL. Specifically, we use the Synthetic Minority Oversampling TEchnique (SMOTE) [27]

and Random Majority Undersampling [26]. SMOTE artificially constructs minority datapoints by interpolating data from known minority cases. In Random Majority Undersampling, majority class instances are randomly selected until a sufficiently large sample is obtained. In addition to these two techniques, we also consider a mixed approach, where we apply half the minority class boost used for full SMOTE, and then apply random majority undersampling to achieve in the end two classes that have approximately equal representation in the training data.

#### D. Classification Algorithms

Finally, we also need to select a machine learning classification algorithm that is able to best predict the validity of a potential link. We consider a wide range of algorithms from families previously used for software engineering applications [28], [29] in order to empirically determine which algorithm is most suitable. Specifically, in this paper we consider k-Nearest Neighbors with  $k = 5$  (5NN), Naïve Bayes (NB), Logistical Regression, Random Forest (RF), Support Vector Machines (SVM), and a Voting ensemble classifier which combines all of the other algorithms. Due to space constraints we do not discuss the implementation of each of these algorithms. Instead we direct the interested reader to a 2010 review of these algorithms [30] for a detailed discussion. Note that for our analysis we use the standard Weka implementation of each algorithm without any special parameter tuning<sup>3</sup>. This allows us to compare classification algorithms in their default state, realizing that additional tuning can further improve performance.

### III. STUDY DESIGN

We performed an empirical study with two main *goals*. The first is to empirically determine the *best configuration for TRAIL* in terms of its ability to automatically support traceability maintenance. The second is to *compare the best configuration of TRAIL to popular IR approaches* that have been previously applied to traceability link recovery and maintenance [5], [6], [7], [17]. In the context of our study, we specifically aim to answer the following research questions:

**RQ<sub>1</sub>:** *What is the best performing configuration for TRAIL?* Given our design space, defined in Table I, we aim to determine the combination of feature selection, rebalancing technique, and classification algorithm that leads to the best performance of TRAIL. We measure performance in terms of F-score, computed as the harmonic mean between precision and recall. We aim to balance recall and precision in order to maximize the number of valid links retrieved by the approach while minimizing the number of false positives contained in the result set.

**RQ<sub>2</sub>:** *Does TRAIL provide superior support for automated traceability maintenance compared with standard IR approaches?* In RQ<sub>1</sub>, we establish a baseline, default configuration for TRAIL chosen to have the highest performance across all datasets under consideration in this study. In RQ<sub>2</sub>, we compare the performance of that default configuration

<sup>3</sup>We used LibSVM as the implementation of the Supported Vector Machine classifier

with IR approaches, which represent the most commonly used techniques for traceability link recovery [4]. In order to provide a conservative analysis of the performance of this default configuration, we directly compare it with the IR technique chosen from the seven approaches in II-A1 to have the highest F-score for each dataset under consideration.

#### A. Data Collection

To answer both research questions, we use a diverse collection of datasets available from the Center of Excellence for Software and Systems Traceability (CoEST) [31] that have been previously used to evaluate new ideas and techniques in the area of traceability [12], [13], [14], [15], [16]. In total, we use 11 datasets that involve eight different types of artifacts, and are extracted from six different software projects. In total, the datasets have 32,616 possible links between pairs of artifacts, of which 2,600 (7.97%) are actually valid and 30,016 are invalid (92.03%). Since the artifacts in these datasets are all either natural language-based (in English or Italian) or source code, we applied the typical preprocessing steps employed in traceability link recovery on this type of artifacts [4]. First, all artifacts were preprocessed to split identifiers on camelCase and underscores ("\_"), then we removed common English (or Italian) words and Java keywords, and finally we stemmed the remaining words to their root form.

Table II depicts the datasets used in our evaluation, containing the number of invalid traceability links, the number of valid traceability links, and the artifact types present in each dataset. Note that, as discussed in section II-C, the data is highly imbalanced, with an average invalid-to-valid link ratio of approximately 12 : 1.

TABLE II  
DATASETS USED IN THE EVALUATION

System	Invalid Links	Valid Links	Artifact Types <sup>†</sup>
eAnci	7091	554 (7.25%)	UC, CC
EasyClinic	1317	93 (6.60%)	UC, CC
EasyClinic	871	69 (7.34%)	ID, CC
EasyClinic	1177	83 (6.59%)	ID, TC
EasyClinic	574	26 (4.33%)	ID, UC
EasyClinic	2757	204 (6.89%)	TC, CC
EasyClinic	1827	63 (3.33%)	TC, UC
eTour	6363	365 (5.43%)	UC, CC
iTrust	1493	58 (3.74%)	UC, CC
MODIS	890	41 (4.40%)	HighR, LowR
SMOS	5656	1044 (15.58%)	UC, CC
<b>Total</b>	<b>30016</b>	<b>2600 (7.97%)</b>	

<sup>†</sup> HighR = High-level Requirements, LowR = Low-level Requirements, UC = Use Cases, CC = Code Classes, ID = Interaction Diagrams, TC = Test Cases.

#### B. Answering RQ<sub>1</sub> - Determining the best configuration of TRAIL

In the first set of experiments, we exhaustively searched the design space for a combination of feature selection, rebalancing technique, and classification algorithm that provides the best performance in terms of F-score across all datasets in the study. We consider this configuration as a baseline TRAIL

TABLE I  
DESIGN SPACE USED FOR THE IMPLEMENTATIONS OF TRAIL

Categories	Variable	Notes
Classification algorithms	5NN	K-nearest neighbors classifier using K=5
	Logistic	Multinomial logistic regression model with a ridge estimator
	NaïveBayes	Naïve Bayes classifier using estimator classes
	Random Forest	Classifier that uses a multitude of random forests
	SVM	Supported Vector Machine classifier
	Voted	Ensembled classifier that averages the output of the previous 5 classifiers to classify links
Feature Selection techniques	none	All attributes are used (i.e., no feature selection is performed)
	cfs	Evaluates the worth of a subset of attributes by considering the individual predictive ability of each feature along with the degree of redundancy between them
	correlation	Evaluates the worth of an attribute by measuring the correlation (Pearson’s) between it and the class
	gainRatio	Evaluates the worth of an attribute by measuring the gain ratio with respect to the class
	infoGain	Evaluates the worth of an attribute by measuring the information gain with respect to the class
	symmetrical	Evaluates the worth of an attribute by measuring the symmetrical uncertainty with respect to the class
Rebalancing techniques	none	No rebalancing technique is applied
	undersampling	Random undersampling is applied to the majority class until it is as small as the minority class
	smote	SMOTE is applied to the minority class until it is as large as the majority class
	5050	First, SMOTE is applied half-way on the minority class, then random undersampling is applied to the majority class until the two classes are equal in size

implementation for  $RQ_2$ . To perform the exhaustive search, we implemented TRAIL using each possible configuration in the design space. We then evaluated each implementation of TRAIL with every dataset independently, running 10-fold cross validations 50 different times and averaging the results. We perform 10-fold cross validations 50 times to account for randomization introduced by the rebalancing techniques and the stratified sampling used in cross validation.

Finally, we collected confusion matrices for each implementation, establishing which configuration achieved the highest average F-score across datasets. We use F-score because it provides a balance between recall (i.e., the ability to extract valid links) and precision (i.e., the ability to minimize false positives). After establishing the best configuration for TRAIL, we used the Mann-Whitney U test (with Holm-Bonferroni correction) to determine whether there is a statistically significant difference between that configuration and the others. We also established effect size using Cliff’s delta.

### C. Answering $RQ_2$ - Comparing TRAIL to IR

In the second set of experiments, we determine the impact of our contribution by comparing the baseline implementation of TRAIL (determined in  $RQ_1$ ) with IR approaches previously used to support TLR. Specifically, we compare against the IR technique among the seven described in Section II-A2 that achieves the highest F-score for each dataset. This leads to a conservative comparison of our baseline configuration.

One important aspect to take into consideration for this comparison is the different nature of the output offered by TRAIL and an IR technique: while the former provides a static set of traceability links predicted as valid, the latter offers a *ranked* list of potential traceability links. To directly compare these two different approaches, we considered each IR technique as a classifier, in which all traceability links in the top  $k$  positions are classified as *valid links*, and those below position  $k$  are classified as *invalid links*. Henceforth, we will refer to the value of  $k$  as the *cut-point*. Note that this

interpretation is specifically for experimental purposes and has previously been used to evaluate similarity-based classifiers in TLR [14]. Since TRAIL predicts a certain number  $N$  of valid traceability links per dataset, we compared TRAIL with each IR technique for every dataset using  $N$  as the cut-point.

Each IR technique provides a single value for each performance metric; for TRAIL, we run 50 trials of 10-fold cross validation on each dataset in order to mitigate the effect of sampling bias and then average the results. We determined if the resulting distributions of precision, recall, and F-score obtained by TRAIL per dataset were normally distributed using the Shapiro-Wilk test with a significance level of 0.01. Next, to compute significance between the distributions provided by TRAIL and the single value of each metric provided by the best IR technique for each dataset, we used either the one-sample t-test (if the observations were normally distributed) or the one-sample Mann-Whitney U test (if the observations were non-normally distributed). Finally, we adjusted the obtained p-values for each metric using the Holm-Bonferroni correction for each dataset. In either case we compute effect sizes by normalizing the difference between the distribution mean and the IR metric by the distribution’s standard deviation. This calculation penalizes non-normally distributed samples, rewarding those with high mean and low standard deviation.

The results of the experiments and the data used in this study are publicly available in a replication package [32].

## IV. RESULTS

### A. $RQ_1$ - Determining the best configuration of TRAIL

Table III shows the average F-score achieved by each combination of parameters in our design space across all of the datasets in our study. The configuration with the highest F-score is displayed in bold with an underline, and is the implementation chosen as a default configuration for use in  $RQ_2$ . The configurations with statistically significantly lower F-scores (at the .01 significance level) compared to the best performing configuration are marked with a superscript symbol,

TABLE III  
AVERAGE F-SCORE (IN PERCENTAGE) ACHIEVED BY THE IMPLEMENTATIONS OF TRAIL ACROSS ALL DATASETS.  
THE BEST CONFIGURATION IS IN UNDERLINED, BOLD FONT.

Rebalancing Technique	Feature Selection	Classifier					
		5NN	Logistic	NaiveBayes	RandomForest	SVM	Vote
none	none	47.43*	50.19*	39.49*	67.18†	0.00*	55.96*
	cfs	59.72*	40.25*	39.22*	63.14*	0.79*	53.84*
	correlation	47.43*	50.19*	39.49*	67.22†	0.00*	56.06*
	gainRatio	61.22*	61.00*	40.17*	72.03	0.00*	66.84*
	infoGain	61.22*	61.00*	40.17*	72.29	0.00*	66.96*
	symmetrical	61.22*	61.00*	40.17*	72.07	0.00*	66.89*
undersampling	none	31.18*	34.60*	36.13*	51.37*	31.22*	38.24*
	cfs	39.88*	37.65*	35.65*	47.43*	34.77*	38.35*
	correlation	31.18*	34.59*	36.13*	51.42*	31.41*	38.28*
	gainRatio	37.63*	38.05*	37.81*	51.38*	35.69*	41.82*
	infoGain	37.63*	38.05*	37.81*	51.34*	35.69*	41.83*
	symmetrical	37.63*	38.05*	37.81*	51.41*	35.69*	41.83*
smote	none	56.19*	56.31*	38.05*	74.80	46.77*	56.94*
	cfs	54.07*	41.04*	37.46*	62.74*	41.44*	45.42*
	correlation	56.15*	56.33*	38.05*	<b>75.18</b>	46.99*	56.99*
	gainRatio	63.10*	58.05*	39.74*	73.89	47.68*	57.74*
	infoGain	63.10*	57.95*	39.75*	73.88	47.67*	57.74*
	symmetrical	63.09*	58.06*	39.77*	73.85	47.69*	57.77*
5050	none	49.59*	51.19*	38.03*	72.09†	43.70*	53.36*
	cfs	51.24*	40.80*	37.35*	61.47*	40.34*	44.33*
	correlation	49.60*	51.17*	38.02*	72.33†	43.97*	53.38*
	gainRatio	57.04*	56.64*	39.67*	70.62†	45.55*	55.01*
	infoGain	57.10*	56.63*	39.69*	70.64†	45.58*	54.99*
	symmetrical	57.04*	56.64*	39.68*	70.60†	45.55*	55.02*

\* = implementations performing statistically significantly worse than the best configuration (0.01 significance level), with a medium or large effect size; † = implementations performing statistically significantly worse than the best configuration (0.01 significance level), with a small effect size.

where an asterisk (“\*”) indicates a medium or large effect size and a cross (“†”) indicates a small effect size. The results indicate that *Random Forrest (RF)* outperforms all other classification algorithms across all other dimensions, and achieves the best results when using the *Pearson correlation* for feature selection and *SMOTE* for data rebalancing.

The findings also show that the choice of *balancing technique* is important, and that undersampling provides the lowest F-scores for all classification algorithms but SVM, across all feature selection techniques. When trying to learn more about this drop in F-score, we noticed that for all configurations not involving SVM, random undersampling boosted recall, but also drastically lowered the precision. Therefore, while undersampling allows TRAIL to retrieve more valid links in those cases, they are returned in addition to many false positives. One reason for this could be that, by reducing the number and diversity of false links to learn from, there were not enough false link instances for TRAIL to learn all the nuances found in the large number of false links in the dataset, therefore misclassifying many of them as true links.

The results also show that completely rebalancing the classes with SMOTE leads to similar F-scores to configurations in which no rebalancing is performed. When analyzing the recall and precision in more detail, we found that SMOTE increases TRAIL’s recall substantially, from 57% (with no rebalancing) to 76% in the case of our selected configuration, while maintaining a precision of 75%. On the other hand, the

selected configuration without any balancing results in a higher precision of 86%, but at 57% recall only slightly more than half of the valid links are extracted. While there are scenarios in which recall should be prioritized over precision or vice versa, in this study we are interested in optimizing the two together. Therefore, practitioners interested in implementing TRAIL should choose an implementation that best suits their needs. For example, in some applications an approach with high recall can be used to reduce the number of document pairs to inspect, while ensuring that a minimal number of valid links are missed. Other applications might call for an approach that maximizes confidence in TRAIL’s predictions by maximizing precision at the cost of overlooking some valid links.

When comparing *feature selection techniques* for RF with SMOTE, which achieves the best results overall, there are small differences in F-score between the parameter choices except for the CFS-Subset algorithm, which achieves a significantly lower F-score. Selection based on Pearson’s correlation provides the best F-score by a narrow margin, and has higher precision than the other options at a slightly reduced recall. There are no statistical differences between configurations using SMOTE with the RF algorithm and any feature selection other than CFS-Subset. Therefore, from a practical point of view, one could also choose the implementation that requires the smallest number of features but statistically provides the same performance. Table IV shows the number of features selected by each of the feature selection techniques for each dataset. As the table

TABLE IV  
NUMBER OF FEATURES SELECTED BY EACH TECHNIQUE PER DATASET. THE INITIAL NUMBER OF FEATURES IS 131

Dataset	cfs	corr	gainR	infoG	symm
eAnci(CC-UC)	8	127	125	125	125
EasyClinic(CC-UC)	10	127	58	58	58
EasyClinic(ID-CC)	6	127	40	40	40
EasyClinic(ID-TC)	9	123	102	102	102
EasyClinic(ID-UC)	9	125	19	19	19
EasyClinic(TC-CC)	10	126	70	70	70
EasyClinic(TC-UC)	6	125	69	69	69
eTour(CC-UC)	10	127	72	72	72
iTrust(CC-UC)	9	126	72	13	13
Modis(HighR-LowR)	3	127	65	65	65
SMOS(CC-UC)	11	125	76	76	76

HighR = High-level Requirements, LowR = Low-level Requirements, UC = Use Cases, CC = Code Classes, ID = Interaction Diagrams, TC = Test Cases.

indicates, the number of features selected is dependent on the system. For this study, however, our goal was to find the best single configuration across all systems. As such, we strictly focus on the implementation that provides the largest average F-score across all systems, i.e., the Pearson’s correlation-based feature selection with full SMOTE and the RF algorithm.

In summary, these findings suggest that the most important consideration in designing a TRAIL implementation is the classification algorithm, where RF shows statistically significantly higher F-scores compared to other algorithms across all other dimensions. Secondly, rebalancing and feature selection are considerations that can be used to arrive at an implementation that favors either recall or precision based on contextual need with the smallest number of features for a given project.

*Correlation-based feature selection with full SMOTE rebalancing and the Random Forest classification algorithm are the best TRAIL configuration for our dataset based on F-score.*

### B. $RQ_2$ - Comparing TRAIL to IR

Table V shows a comparison of precision, recall, and F-score between our derived baseline TRAIL implementation and the best performing IR technique for each dataset from the set of seven techniques considered in our study. Again, the best performing IR approach was chosen by selecting the one with the highest F-score, as done for the default configuration of TRAIL. For IR, the table shows precision, recall and F-score at K, where K is the cut-point chosen to be the number of links TRAIL predicts to be valid (which allows for a direct comparison between the approaches, as described in section III-C). The IR techniques that performed statistically significantly worse than TRAIL with a significance level of 0.01 and large effect sizes or greater are denoted with an asterisk (\*). Interestingly, VSM, one of the most basic IR techniques considered in this study, had the highest F-score in nine of twelve systems, whereas LDA and LM-JM did not have the highest F-score for any of the cases. This is consistent with previous work on IR-based TLR [17].

Our baseline TRAIL implementation outperforms even the best IR approach in terms of precision, recall, and F-score for

each of the 11 datasets in this study. However, we note that the performance of both TRAIL and IR is dependent on the dataset. For six datasets, we achieve higher than 70% recall and precision, while the lowest recall and precision for TRAIL is 56% and 59% respectively for EasyClinic(ID-UC). For IR, recall and precision higher than 60% is only achieved in one case (EasyClinic(CC-UC)), and the lowest recall and precision are 30% and 28% respectively, more than 25 percentage points lower than that of TRAIL.

There is some variability in the performance improvement provided by TRAIL. For example, in the case of EasyClinic(CC-UC) and iTrust, TRAIL outperforms VSM by less than five percentage points. However, on the other end of the spectrum, TRAIL is able to extract more than 90% of the valid links (with a maximum of 99%) for three of the EasyClinic datasets with higher than 90% precision. This is a substantial improvement over IR approaches, which never achieve more than 60% recall or precision for those datasets. Similarly, TRAIL improves performance for eAnci in terms of each metric by more than 40 percentage points. Further, in six of the datasets, TRAIL outperforms IR by more than 30 percentage points for all three metrics. Overall, our baseline TRAIL configuration outperforms IR by more than 26 percentage points in average precision, recall, and F-score. Finally, the performance of TRAIL is statistically significantly better than IR at the .01 significance level for all three metrics for each dataset, with a large effect size or greater in each case. Therefore, while many of the features used by our baseline implementation of TRAIL are derived directly from IR techniques, our findings indicate that the power of TRAIL may come from combining this type of information with the other features used.

In summary, the results of our study indicate that even our baseline implementation of TRAIL outperforms traditional IR approaches for all 11 datasets in terms of all three performance metrics. Further, while the improvement in performance is limited in some cases, in the majority of cases TRAIL improves each measure by more than 10 percentage points. Moreover, these findings are statistically significant at the .01 significance level, with a large effect size or greater. Considering all datasets in the study, TRAIL outperforms even the best IR for each dataset by more than 26 percentage points in terms of average precision, recall, and F-score. Given that we did not optimize TRAIL per dataset, but rather used a baseline, general implementation across all systems, we anticipate that further refining our parameters per dataset will lead to even better results for TRAIL in future work.

TRAIL significantly outperforms the best Information Retrieval approaches in terms of precision, recall, and F-score in all datasets.

## V. THREATS TO VALIDITY

*Construct validity* refers to how well the metrics used for evaluation truly capture what they were intended to. We mitigated this risk in several different ways. First, we measured performance of both TRAIL and IR techniques throughout



TABLE V  
PRECISION, RECALL AND F-SCORE (IN PERCENTAGE) FOR TRAIL COMPARED TO THE IR TECHNIQUE WITH THE HIGHEST F-SCORE FOR EACH DATASET.

Dataset	TRAIL			Information Retrieval			
	Precision(%)	Recall(%)	F-score(%)	Best IR	Precision@K(%)	Recall@K(%)	F-score@K(%)
eAnci(CC-UC)	75.46	80.53	77.91	VSM	28.21*	30.14*	29.14*
EasyClinic(CC-UC)	64.06	71.46	67.54	VSM	61.54*	68.82*	64.97*
EasyClinic(ID-CC)	72.56	73.19	72.83	VSM	58.57*	59.42*	58.99*
EasyClinic(ID-TC)	94.14	92.75	93.43	LSA	58.54*	57.83*	58.18*
EasyClinic(ID-UC)	59.27	56.23	57.59	VSM	48.00*	46.15*	47.06*
EasyClinic(TC-CC)	95.96	99.03	97.47	VSM	44.08*	45.59*	44.82*
EasyClinic(TC-UC)	93.39	95.62	94.47	LM-Dirichlet	55.38*	57.14*	56.25*
eTour(CC-UC)	57.17	64.98	60.82	VSM	49.40*	56.16*	52.56*
iTrust(CC-UC)	56.79	65.76	60.92	VSM	52.94*	62.07*	57.14*
Modis(HighR-LowR)	65.88	62.93	64.32	VSM	32.50*	31.71*	32.10*
SMOS(CC-UC)	87.07	73.53	79.73	BM25	37.76*	31.90*	34.58*
<b>Average</b>	<b>74.70</b>	<b>76.00</b>	<b>75.18</b>		<b>47.90</b>	<b>49.72</b>	<b>48.71</b>

\* = statistical significance (0.01 significance level) with a large effect size or greater.

the analysis using metrics that have been commonly used in software engineering research, and specifically in TLR studies [33], [4]. Second, for both RQ<sub>1</sub> and RQ<sub>2</sub>, we considered the effect of randomization on the performance metrics as it is possible using only one trial to obtain good results purely by chance. Therefore, we performed 50 trials of TRAIL for each experiment, averaging to aggregate our final results. Finally, we employ statistical tests where applicable, adding rigor to the analysis and empirically supporting our claims.

*Internal validity* refers to how well a study mitigates multiple independent (i.e., conflating) variables from interfering with inferences on the dependent variable. We mitigate these risks in RQ<sub>1</sub> by exhaustively searching the design space, systematically varying independent variables one at a time. Further, in answering this research question we derive a baseline implementation of TRAIL which is then compared to the best performing IR approach for each dataset. Finally, in RQ<sub>2</sub>, we carefully constructed a methodology that allows direct comparisons between TRAIL and IR despite their fundamentally different approaches to extracting traceability links and corresponding output.

*External validity* investigates how well the findings of a study can be generalized. For this study we mitigated this risk by considering a diverse set of 11 datasets extracted from six different software systems and eight types of artifacts and unique artifact combinations, involving more than 30,000 potential traceability links. Moreover, these are all datasets curated by the research community, which have been used in previous traceability studies.

## VI. RELATED WORK

Broadly speaking, software traceability covers a wide range of contexts linking abstract concepts such as architectural tactics and non-functional requirements (NFRs) [34], [35] or software artifacts like use cases and source code [36], [37]. Significant previous work has been done to increase the level of automation available for establishing and maintaining these links, most commonly using IR techniques to rank artifacts or candidate links based on document similarity. Here we discuss four main

categories of approaches: IR, machine learning or classification, event-driven, and model-based approaches.

Applications of IR to TLR are the most related to this work, and began with probabilistic [36], [38] and vector space models [39] to retrieve links between source code and documentation, as well as source code and requirements. Additionally, other IR approaches such as Latent Semantic Analysis (LSA) [5], probabilistic LSA (pLSA) [6], Jensen-Shannon (JS)[6], Latent Dirichlet Allocation (LDA) [7], and proximity-based VSM [8] have also been applied directly to the TLR task. Oliveto et al. performed an empirical study of IR methods for TLR [17] comparing VSM, LSA, LDA, and JS via Principal Component Analysis, showing that VSM, LSA, and JS capture similar information, while information captured by LDA is unique. Learning to Rank has been applied to improve IR approaches [37], deriving a more accurate ranking of document similarities by combining information from various IR approaches into a single list. More recently, Falessi et al. [40] used machine learning classifiers to estimate the number of valid links remaining in a set of candidate links returned by IR techniques.

There are also a large number of techniques for increasing the performance of IR approaches to TLR by augmenting document information with key-phrases and technical terms through a thesaurus [12], applying alternative term weighting strategies [41], [42], and smoothing filters [43]. Additionally, Cleland-Huang et al. [44] propose three different improvement strategies based on artifact hierarchies, clustering, and graph pruning. Further, previous research has considered the selection of specific IR-based infrastructures for TLR through the use of genetic algorithms [45]. For a more in-depth overview of related work on IR in TLR, we direct the reader to a systematic mapping study of the field [4]. Our work differs from this previous work as we leverage IR techniques to generate features for our classifier rather than applying them directly to the task of establishing traceability links.

There is also existing work specifically focused on classification and traceability, but significantly different than our approach. In the area of requirements engineering, Cleland-Huang et al. proposed a probabilistic classifier trained on a

set of indicator words for NFRs [34], which was shown to outperform traditional machine learning classifiers [46]. This classifier was applied to traceability specifically for linking regulatory codes to project-specific requirements [14] and architectural tactics to source code [35]. A semi-supervised, iterative approach utilizing Expectation Maximization has also been proposed [47]. Further, Asuncion et al. constructed an automatic approach called TRASE using probabilistic topic modeling with a modified form of LDA to automatically infer traceability links [7].

Our work differs from these classification schemes in a couple of ways. The approach by Cleland-Huang et al. [14] requires the identification of a set of high level concepts that can be represented by an indicator term set. Our approach does not require this intermediate step and is able to directly classify links between artifacts. Furthermore, their approach uses a threshold for classification, similar to IR-based approaches, while our classification criteria are left entirely to the machine learning algorithm, making it fully automatic. Our approach differs from that by Casamayor et al. [47] in that they provide an iterative approach based on user feedback for requirement classification, while we propose a customizable, yet automatic approach for the classification of traceability links, and in a more generalized context. Compared to TRASE [7], TRAIL uses LDA as a single IR feature (along many others) to determine document similarity, which is then used by the larger classification model to learn relationships between artifacts.

While the aforementioned approaches can be applied to the maintenance problem by using them to re-establish traceability after changes are made, there have also been event-driven techniques specifically for maintaining existing traceability links. Trigger events for re-evaluating links during system evolution were suggested by Chen and Chou as early as 1999 [48]. Cleland-Huang et al. [49] followed with an approach based on an event server intended to notify stakeholders when related artifacts are updated and traceability information might be stale. Murta et al. [50] use a rule-based approach to ensure that traceability links are consistent throughout the evolution process specifically for maintaining links between architectural concepts and implementing code. Mäder et al. [51] provide a modeling plugin that detects high-level operations on software models and propagates appropriate updates to underlying traceability links with related artifacts. Similarly, TraceMAINTAINER [52], [53] is a semi-automated maintenance system that also uses events and a rules engine to perform updates. While these approaches specifically seek to address the maintenance problem, they have some significant limitations. First, they typically require a set of rules to be established by which updates are propagated when an event is raised. Second, these approaches often seek to prevent trace decay in which existing links become stale. Therefore, they do not address the situation in which new artifacts are created and must be included in an existing web of traceability links. Our work differs from these approaches in that we use historical training data to maintain links rather than relying on predefined rule sets. Additionally, our approach supports both the situation in which existing links are altered

and when new links are established as artifacts are created.

Finally, there has been a significant amount of work in model-driven engineering around establishing and maintaining traceability links between models. These approaches are related to this work only in that they seek to address the same general problem, but with techniques that are not typically applicable outside of model-driven projects, while TRAIL is fully generalizable. As such, we refer the interested reader to any of several surveys on related topics [54], [55], [56].

## VII. ACKNOWLEDGEMENTS

This work is supported in part by the National Science Foundation (NSF) grant CCF-1526929.

## VIII. CONCLUSION AND FUTURE WORK

In this paper, we propose TRAIL, a novel technique for automating traceability maintenance by considering TLR as a binary classification problem. We address this problem using machine learning algorithms trained on historical traceability information. We empirically derive a high-performance baseline implementation of TRAIL on 11 datasets commonly used to evaluate new approaches to traceability. Moreover, we show that TRAIL significantly outperforms even the best IR approach for each of these datasets as measured by three performance metrics: precision, recall, and F-score.

While a significant improvement over existing IR techniques, TRAIL suffers from some of the same limitations as traditional IR approaches to traceability. For example, this study did not address the so-called *vocabulary mismatch* problem, which arises when different terms are used to express the same ideas between sets of artifacts. Recent work in traceability [57] has presented techniques for bridging the term gap between artifact sets. Our future work will focus on improving the IR components of our feature set by incorporating state-of-the-art improvements such as this.

Moreover, because at its core TRAIL is a direct application of machine learning classification to the TLR problem, we can apply the most recent advances in this field to expand the framework's capabilities. First, we plan to investigate transfer-learning techniques that make it possible to train a model using historical data from another, similar project. The result is an implementation of TRAIL more widely applicable to greenfield projects with no existing historical data. Further, we will investigate active learning approaches that look to minimize the amount of training data required to generate an effective predictive model for TRAIL. Using these techniques in tandem with transfer learning, we can minimize barriers to adopting TRAIL from a practitioner's perspective. Finally, the existing framework is based on traditional machine learning, which requires feature engineering to derive vector representations of the traceability links. Future work will investigate using deep neural networks capable of automatically extracting vector representations of links for classification, therefore eliminating the need to design features.

## REFERENCES

- [1] E. Bouillon, P. Mäder, and I. Philippow, "A survey on usage scenarios for requirements traceability in practice," in *International Working Conference on Requirements Engineering: Foundation for Software Quality*. Springer, 2013, pp. 158–173.
- [2] P. Mäder and A. Egyed, "Do developers benefit from requirements traceability when evolving and maintaining a software system?" *Empirical Software Engineering*, vol. 20, no. 2, pp. 413–441, 2015.
- [3] P. Rempel and P. Mader, "Preventing defects: The impact of requirements traceability completeness on software quality," *IEEE Transactions on Software Engineering*, 2016.
- [4] M. Borg, P. Runeson, and A. Ardö, "Recovering from a decade: A systematic mapping of information retrieval approaches to software traceability," *Empirical Software Engineering*, vol. 19, no. 6, pp. 1565–1616, 2014.
- [5] A. Marcus, J. I. Maletic, and A. Sergeyeu, "Recovery of traceability links between software documentation and source code," *International Journal of Software Engineering and Knowledge Engineering*, vol. 15, no. 05, pp. 811–836, 2005.
- [6] A. Abadi, M. Nisenson, and Y. Simionovici, "A traceability technique for specifications," in *Proceedings of the 16th IEEE International Conference on Program Comprehension (ICPC'08)*, 2008, pp. 103–112.
- [7] H. U. Asuncion, A. U. Asuncion, and R. N. Taylor, "Software traceability with topic modeling," in *Proceedings of the 32nd International Conference on Software Engineering*, vol. 1. IEEE, 2010, pp. 95–104.
- [8] W. Kong and J. H. Hayes, "Proximity-based traceability: An empirical validation using ranked retrieval and set-based measures," in *Proceedings of the 1st International Workshop on Empirical Requirements Engineering (EmpiRE)*, 2011, pp. 45–52.
- [9] L. James, "Automatic requirements specification update processing from a requirements management tool perspective," in *Proceedings of the International Conference and Workshop on Engineering of Computer-Based Systems*. IEEE, 1997, pp. 2–9.
- [10] K. Weidenhaupt, K. Pohl, M. Jarke, and P. Haumer, "Scenarios in system development: current practice," *IEEE Software*, vol. 15, no. 2, pp. 34–45, 1998.
- [11] G. Antoniol, C. Casazza, and A. Cimitile, "Traceability recovery by modeling programmer behavior," in *Proceedings of the Seventh Working Conference on Reverse Engineering*. IEEE, 2000, pp. 240–247.
- [12] J. H. Hayes, A. Dekhtyar, and J. Osborne, "Improving requirements tracing via information retrieval," in *Proceedings of the 11th IEEE International Requirements Engineering Conference (RE'03)*, 2003, pp. 138–147.
- [13] J. H. Hayes, A. Dekhtyar, and S. K. Sundaram, "Advancing candidate link generation for requirements tracing: The study of methods," *Transactions Software Engineering*, vol. 32, no. 1, pp. 4–19, 2006. [Online]. Available: <http://dx.doi.org/10.1109/TSE.2006.3>
- [14] J. Cleland-Huang, A. Czauderna, M. Gibiec, and J. Emenecker, "A machine learning approach for tracing regulatory codes to product specific requirements," in *Proceedings of the International Conference on Software Engineering*, vol. 1.
- [15] N. Ali, Y.-G. Gueheneuc, and G. Antoniol, "Requirements traceability for object oriented systems by partitioning source code," in *Proceedings of the 18th Working Conference on Reverse Engineering (WCRE'11)*, 2011, pp. 45–54.
- [16] M. Gethers, R. Oliveto, D. Poshvanyk, and A. De Lucia, "On integrating orthogonal information retrieval methods to improve traceability recovery," in *IEEE International Conference on Software Maintenance (ICSM)*. IEEE, 2011, pp. 133–142.
- [17] R. Oliveto, M. Gethers, D. Poshvanyk, and A. De Lucia, "On the equivalence of information retrieval methods for automated traceability link recovery," in *IEEE International Conference on Program Comprehension (ICPC)*, 2010, pp. 68–71.
- [18] C. Mills and S. Haiduc, "The impact of retrieval direction on ir-based traceability link recovery," in *39th International Conference on Software Engineering: New Ideas and Emerging Technologies Results Track*. IEEE, 2017, pp. 51–54.
- [19] C. Zhai and J. Lafferty, "A study of smoothing methods for language models applied to ad hoc information retrieval," in *Proceedings of the 24th annual international ACM SIGIR Conference on Research and Development in Information Retrieval*, 2001, pp. 334–342.
- [20] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman, "Indexing by Latent Semantic Analysis," *Journal of the American Society for Information Science*, vol. 41, no. 6, pp. 391–407, 1990.
- [21] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *The Journal of machine Learning research*, vol. 3, pp. 993–1022, 2003. [Online]. Available: <http://dl.acm.org/citation.cfm?id=944937>
- [22] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [23] S. Haiduc, G. Bavota, R. Oliveto, A. De Lucia, and A. Marcus, "Automatic query performance assessment during the retrieval of software artifacts," in *Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering (ASE'12)*, 2012, pp. 90–99.
- [24] S. Haiduc, G. De Rosa, G. Bavota, A. Marcus, R. Oliveto, and A. De Lucia, "Query quality prediction and reformulation for source code search: the Refoqus tool," in *Proceedings of the 35th IEEE/ACM International Conference on Software Engineering (ICSE'13)*, San Francisco, USA, 2013, pp. 1307–1310.
- [25] C. Mills, G. Bavota, S. Haiduc, R. Oliveto, A. Marcus, and A. D. Lucia, "Predicting query quality for applications of text retrieval to software engineering tasks," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 26, no. 1, p. 3, 2017.
- [26] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 21, no. 9, pp. 1263–1284, Sep. 2009. [Online]. Available: <http://dx.doi.org/10.1109/TKDE.2008.239>
- [27] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: Synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [28] E. Giger, M. D'Ambros, M. Pinzger, and H. C. Gall, "Method-level bug prediction," in *Proceedings of the ACM-IEEE international symposium on Empirical software engineering and measurement*. ACM, 2012, pp. 171–180.
- [29] C. Seiffert, T. M. Khoshgoftaar, J. Van Hulse, and A. Folleco, "An empirical study of the classification performance of learners on imbalanced and noisy software quality data," *Information Sciences*, vol. 259, pp. 571–595, 2014.
- [30] A. Khan, B. Baharudin, L. H. Lee, and K. Khan, "A review of machine learning algorithms for text-documents classification," *Journal of Advances in Information Technology*, vol. 1, no. 1, pp. 4–20, 2010.
- [31] *COEST: Community Data Sets*, 2017. [Online]. Available: <http://coest.org>
- [32] *Replication Package*, 2018. [Online]. Available: [http://www.cs.fsu.edu/~serene/mills\\_icsme\\_18\\_trace/](http://www.cs.fsu.edu/~serene/mills_icsme_18_trace/)
- [33] J. Cleland-Huang, B. Berenbach, and S. Clark, "Best Practices for Automated Traceability," *Computer*, no. 40(6), pp. 27–35, jun 2007.
- [34] J. Cleland-Huang, R. Settini, X. Zou, and P. Solc, "Automated classification of non-functional requirements," *Requirements Engineering*, vol. 12, no. 2, pp. 103–120, 2007.
- [35] M. Mirakhorli, Y. Shin, J. Cleland-Huang, and M. Cinar, "A tactic-centric approach for automating traceability of quality concerns," in *Proceedings of the International Conference on Software Engineering (ICSE)*, 2012, pp. 639–649.
- [36] G. Antoniol, G. Canfora, G. Casazza, A. De Lucia, and E. Merlo, "Tracing object-oriented code into functional requirements," in *Proceedings of the International Workshop on Program Comprehension*, 2000, pp. 79–86.
- [37] D. Binkley and D. Lawrie, "Learning to rank improves IR in SE," in *Proceedings of the 30th IEEE International Conference on Software Maintenance and Evolution (ICSME'14)*. IEEE, 2014, pp. 441–445.
- [38] G. Antoniol, G. Canfora, A. De Lucia, and E. Merlo, "Recovering code to documentation links in oo systems," in *Proceedings of the Sixth Working Conference on Reverse Engineering*, 1999, pp. 136–144.
- [39] G. Antoniol, G. Canfora, G. Casazza, A. De Lucia, and E. Merlo, "Recovering traceability links between code and documentation," *IEEE Transactions Software Engineering*, vol. 28, no. 10, pp. 970–983, 2002.
- [40] D. Falessi, M. Di Penta, G. Canfora, and G. Cantone, "Estimating the number of remaining links in traceability recovery," *Empirical Software Engineering*, pp. 1–32, 2016.
- [41] R. Settini, J. Cleland-Huang, O. B. Khadra, J. Mody, W. Lukasik, and C. DePalma, "Supporting software evolution through dynamically retrieving traces to UML artifacts," in *Proceedings of the 7th International Workshop on Principles of Software Evolution*, 2004, pp. 49–54.
- [42] S. K. Sundaram, J. H. Hayes, A. Dekhtyar, and E. A. Holbrook, "Assessing traceability of software engineering artifacts," *Requirements Engineering*, vol. 15, no. 3, pp. 313–335, 2010.

- [43] A. De Lucia, M. Di Penta, R. Oliveto, A. Panichella, and S. Panichella, "Improving IR-based traceability recovery using smoothing filters," in *Proceedings of the 19th IEEE International Conference on Program Comprehension (ICPC'11)*, 2011, pp. 21–30.
- [44] J. Cleland-Huang, R. Settimi, C. Duan, and X. Zou, "Utilizing supporting evidence to improve dynamic requirements traceability," in *Proceedings of the 13th IEEE International Requirements Engineering Conference (RE'05)*, 2005, pp. 135–144.
- [45] S. Lohar, S. Amornborvornwong, A. Zisman, and J. Cleland-Huang, "Improving trace accuracy through data-driven configuration and composition of tracing features," in *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*. ACM, 2013, pp. 378–388.
- [46] A. Jalaji, R. Goff, M. Jackson, N. Jones, and T. Menzies, "Making sense of text: identifying non functional requirements early," *West Virginia University CSEE Technical Report*, 2006.
- [47] A. Casamayor, D. Godoy, and M. Campo, "Identification of non-functional requirements in textual specifications: A semi-supervised learning approach," *Information and Software Technology*, vol. 52, no. 4, pp. 436–445, 2010.
- [48] J.-Y. Chen and S.-C. Chou, "Consistency management in a process environment," *Journal of Systems and Software*, vol. 47, no. 2-3, pp. 105–110, 1999.
- [49] J. Cleland-Huang, C. K. Chang, and J. C. Wise, "Automating performance-related impact analysis through event based traceability," *Requirements Engineering*, vol. 8, no. 3, pp. 171–182, 2003.
- [50] L. G. Murta, A. van der Hoek, and C. M. Werner, "Continuous and automated evolution of architecture-to-implementation traceability links," *Automated Software Engineering*, vol. 15, no. 1, pp. 75–107, 2008.
- [51] P. Mäder, O. Gotel, and I. Philippow, "Rule-based maintenance of post-requirements traceability relations," in *International Requirements Engineering*. IEEE, 2008, pp. 23–32.
- [52] P. Mader, O. Gotel, and I. Philippow, "Enabling automated traceability maintenance by recognizing development activities applied to models," in *Proceedings of the 23rd IEEE/ACM International Conference on Automated Software Engineering*. IEEE Computer Society, 2008, pp. 49–58.
- [53] P. Mäder, O. Gotel, and I. Philippow, "Enabling automated traceability maintenance through the upkeep of traceability relations," in *European Conference on Model Driven Architecture-Foundations and Applications*. Springer, 2009, pp. 174–189.
- [54] I. Galvao and A. Goknil, "Survey of traceability approaches in model-driven engineering," in *Enterprise Distributed Object Computing Conference, 2007. EDOC 2007. 11th IEEE International*. IEEE, 2007, pp. 313–313.
- [55] S. Winkler and J. von Pilgrim, "A survey of traceability in requirements engineering and model-driven development," *Software & Systems Modeling*, vol. 9, no. 4, pp. 529–565, 2010.
- [56] I. Boussaïd, P. Siarry, and M. Ahmed-Nacer, "A survey on search-based model-driven engineering," *Automated Software Engineering*, vol. 24, no. 2, pp. 233–294, 2017.
- [57] J. Guo, M. Gibiec, and J. Cleland-Huang, "Tackling the term-mismatch problem in automated trace retrieval," *Empirical Software Engineering*, vol. 22, no. 3, pp. 1103–1142, 2017.