



Traceability Link Recovery between Requirements and Models using an Evolutionary Algorithm Guided by a Learning to Rank Algorithm: Train control and management case

Ana C. Marcén^{a,b,*}, Raúl Lapeña^a, Óscar Pastor^b, Carlos Cetina^a

^aSVIT Research Group, Universidad San Jorge, Spain

^bCentro de Investigación en Métodos de Producción de Software, Universitat Politècnica de València, Spain

ARTICLE INFO

Article history:

Received 9 March 2018

Revised 30 September 2019

Accepted 10 January 2020

Available online 15 January 2020

Keywords:

Traceability Link Recovery

Requirements Engineering

Models

Evolutionary Algorithm

Learning to Rank

ABSTRACT

Traceability Link Recovery (TLR) has been a topic of interest for many years within the software engineering community. In recent years, TLR has been attracting more attention, becoming the subject of both fundamental and applied research. However, there still exists a large gap between the actual needs of industry on one hand and the solutions published through academic research on the other.

In this work, we propose a novel approach, named Evolutionary Learning to Rank for Traceability Link Recovery (TLR-ELtoR). TLR-ELtoR recovers traceability links between a requirement and a model through the combination of evolutionary computation and machine learning techniques, generating as a result a ranking of model fragments that can realize the requirement.

TLR-ELtoR was evaluated in a real-world case study in the railway domain, comparing its outcomes with five TLR approaches (Information Retrieval, Linguistic Rule-based, Feedforward Neural Network, Recurrent Neural Network, and Learning to Rank). The results show that TLR-ELtoR achieved the best results for most performance indicators, providing a mean precision value of 59.91%, a recall value of 78.95%, a combined F-measure of 62.50%, and a MCC value of 0.64. The statistical analysis of the results assesses the magnitude of the improvement, and the discussion presents why TLR-ELtoR achieves better results than the baselines.

© 2020 Elsevier Inc. All rights reserved.

1. Introduction

Traceability Link Recovery (TLR) has been a subject of investigation for many years within the software engineering community (Gotel and Finkelstein, 1994; Spanoudakis and Zisman, 2005). Research has shown that affordable traceability can be critical to the success of a project (Watkins and Neal, 1994) and leads to increased maintainability and reliability of software systems by making it possible to verify and to trace non-reliable parts (Ghazarian, 2010). Specifically, more complete traceability decreases the expected defect rate in the developed software (Rempel and Mäder, 2017).

In recent years, TLR has been attracting more attention, and re-establishing the traceability links between software artifacts has become a subject of both fundamental and applied research (Parizi et al., 2014). In fact, a few approaches have been proposed to recover traceability between requirements and models

(Winkler and Pilgrim, 2010). However, the support of traceability research for practical problems in industry is perceived as being rather low (Winkler and Pilgrim, 2010), and there still exists a large gap between the needs of industry on one hand and the published solutions from academic research on the other (Rummeler et al., 2007).

In this work, we propose a novel approach, named Evolutionary Learning to Rank for Traceability Link Recovery (TLR-ELtoR). TLR-ELtoR recovers traceability links between the requirements of a software system and the models that implement it. Specifically, our approach is based on an Evolutionary Algorithm (EA). Moreover, the EA is guided by a Learning to Rank algorithm that empowers us to take advantage of the knowledge and the experience that have been generated in companies for years in order to automatically perform ranking tasks. In summary, from a requirement and a model, TLR-ELtoR generates a ranking of model fragments that can realize the requirement.

* Corresponding author.

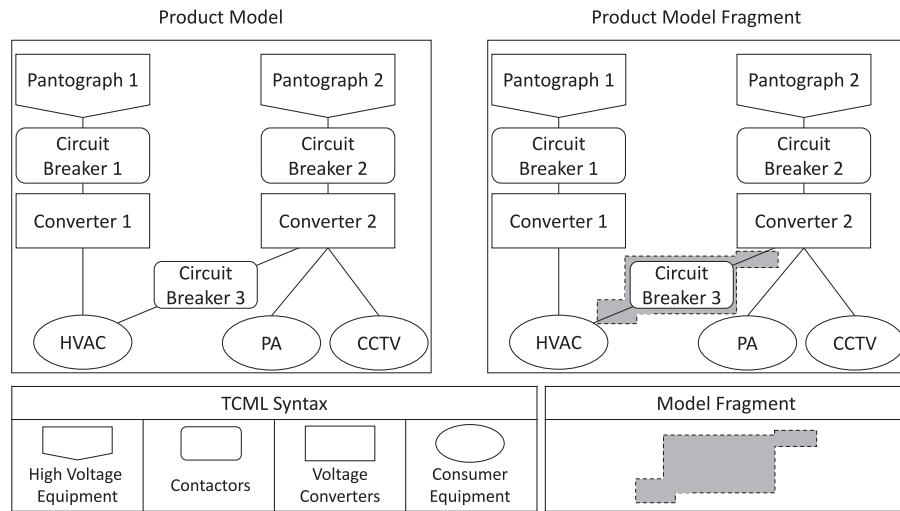


Fig. 1. Example of a TCML model and model fragment

The presented approach was evaluated in a real-world case study provided by our industrial partner, CAF¹ (Construcciones y Auxiliar de Ferrocarriles), a worldwide provider of railway solutions. The outcomes of TLR-ELtoR were compared with five TLR approaches, these approaches were selected taking into account the approaches that obtain the best results for recovering traceability between requirements and models Winkler and Pilgrim (2010), the successful application of deep learning techniques for TLR in recent works Guo et al. (2017), and the size of the search space to be explored. The first one (Spanoudakis et al., 2004) is a Linguistic Rule-Based (TLR-Linguistic) approach that is based on Parts-of-Speech (POS) Tagging and traceability rules. The second one (De Lucia et al., 2004; Lucia et al., 2007) is an Information Retrieval (TLR-IR) approach that is based on Latent Semantic Indexing (LSI) and Singular Value Decomposition (SVD). The third one is a Feedforward Neural Network (TLR-FNN) approach that is based on a traditional neural network structure. The fourth one is a Recurrent Neural Network (TLR-RNN) approach that is based on an extension of a Feedforward Neural Network with feedback connections to model the temporal characteristics of the problem being learned (Engelbrecht, 2007). The fifth one is a Learning to Rank (TLR-LtoR) approach based on ranking Machine Learning algorithms of the same name.

The results show that TLR-ELtoR achieved the best results for most the performance indicators, providing a mean precision value of 59.91%, a recall value of 78.95%, a combined F-measure of 62.50%, and a MCC value of 0.64. In contrast, the TLR-Linguistic baseline, the TLR-IR baseline, and the TLR-LtoR baseline had worse results for these same measurements. On the other hand, although TLR-FNN and TLR-RNN achieved the best results for recall, they obtained the worst results for the rest of the indicators. The statistical analysis of the results assesses the magnitude of the improvement, and the discussion presents two advantages and a limitation of our approach.

The remainder of this paper is structured as follows: Section 2 provides background on our case study. Section 3 highlights our TLR-ELtoR approach. Sections 4 and 5 detail the genetic operations step and the fitness function step of our approach, respectively. Section 6 details the means used to evaluate our work and the results of the evaluation. Section 7 analyzes the statistical significance of the obtained results. Section 8 discusses our approach and the obtained results. Section 9 describes the threats to

the validity of our work. Section 10 introduces the existing works that are related to our work. Section 11 provides the means to replicate the results. Finally, Section 12 concludes the paper.

2. Background

This section presents the Train Control and Management Language (TCML), which is used to formalize the products manufactured by our industrial partner. TCML has the expressiveness required to describe the interaction between the main pieces of equipment installed in a train unit. TCML also has the required expressiveness to specify non-functional aspects that are related to regulation, such as the quality of signals from the equipment or the different levels of installed redundancy. TCML will be used through the rest of the paper to present a running example. In this work, for the sake of the understandability and legibility of the running example, we present an equipment-focused, simplified subset of TCML, with four different kinds of equipment:

1. **High voltage equipment**, which is in charge of harvesting the energy that powers the different elements of the train.
2. **Contactors**, which are in charge of opening or closing the circuits between the High Voltage equipment and the Voltage Converters.
3. **Voltage converters**, which are in charge of transforming the harvested electric power into a current that the Consumer Equipment can work with.
4. **Consumer equipment**, which is in charge of carrying out all of the tasks required for the train to work properly and provide comfort to the passengers.

Fig. 1 depicts an example that is taken from a real-world train. It presents a converter assistance scenario. In the example, two separate pantographs (High Voltage Equipment) collect energy from the overhead wires and send it to their respective circuit breakers (Contactors), which in turn send it to their independent Voltage Converters. The converters then power their assigned Consumer Equipment: the left one powers the HVAC (the air conditioning system of the train) devices, and the right one powers the PA (public address system) and the CCTV (television system) circuits of the train.

There is an additional circuit breaker between the second converter and the HVAC that is connected to the first converter. The part on the right of Fig. 1 shows an example of a model fragment

¹ www.caf.net/en.

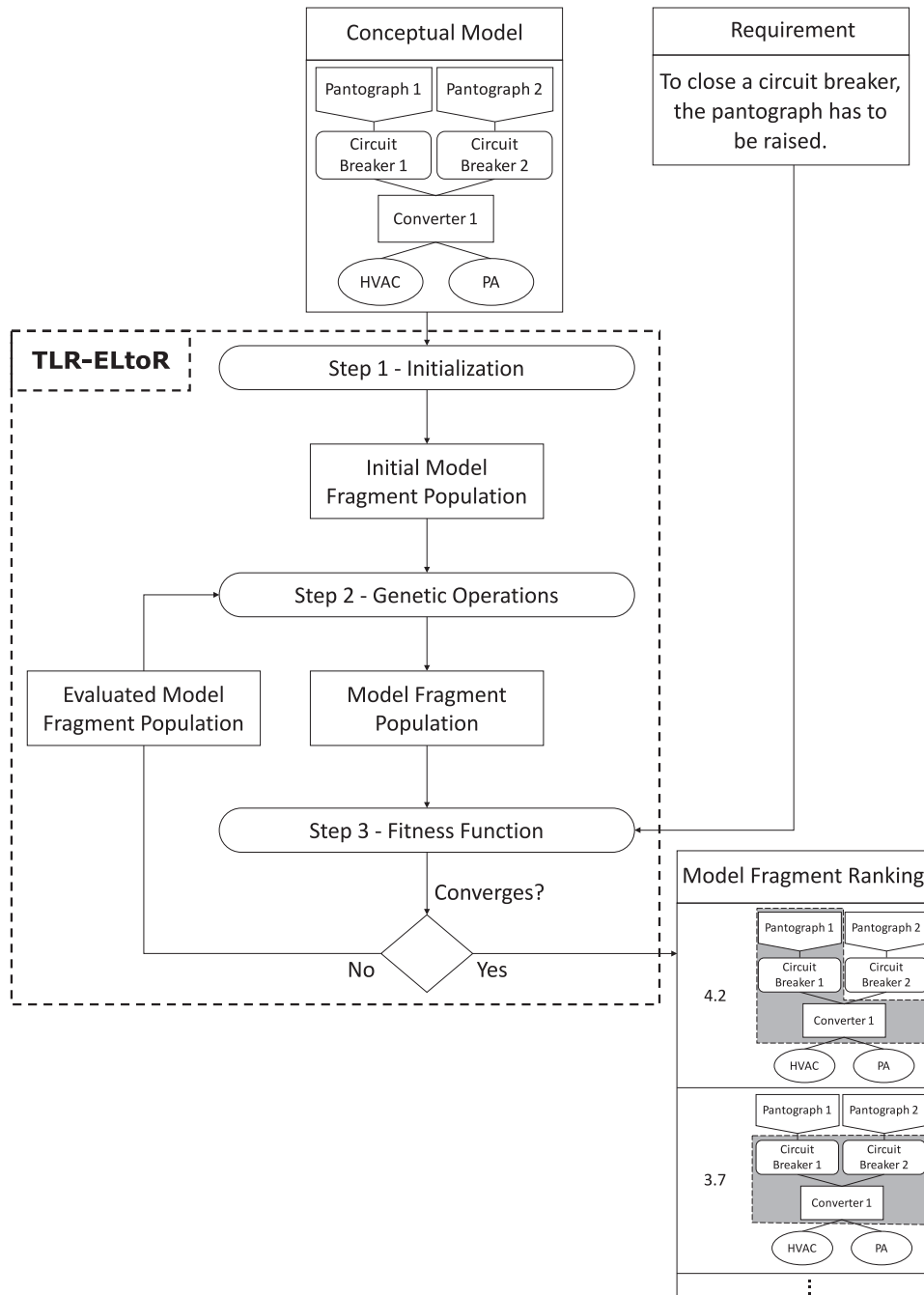


Fig. 2. Overview of the approach

of the product model. The model fragment includes the additional circuit breaker. This model fragment is the realization of the converter assistance requirement, which allows the passing of current from one converter to a piece of Consumer Equipment that is assigned to its peer. In the case of overload or failure of the first converter, total or partial functioning (depending on specific conditions) of the HVAC could be covered by the second converter.

To formalize the model fragments, we use the Common Variability Language (CVL) (Haugen et al., 2008). CVL defines variants of a base model (conforming to MOF, the OMG metalanguage for defining modeling languages) by replacing variable parts of the base model with alternative model replacements that are found in a library.

3. Overview of our TLR-ELtoR approach

This section presents the proposed TLR-ELtoR approach for TLR between the requirements and the models through an EA, which is based on genetic operations and a fitness function. The objective of the approach is to provide the model fragment from a given model that realizes a specific requirement. To do this, the approach receives as input the model that implements a specific requirement. The approach relies on an evolutionary algorithm that iterates over a population of model fragments, evolving them using genetic operations. Then, the score of each model fragment and its position in the ranking are calculated through the fitness function that uses Learning to Rank

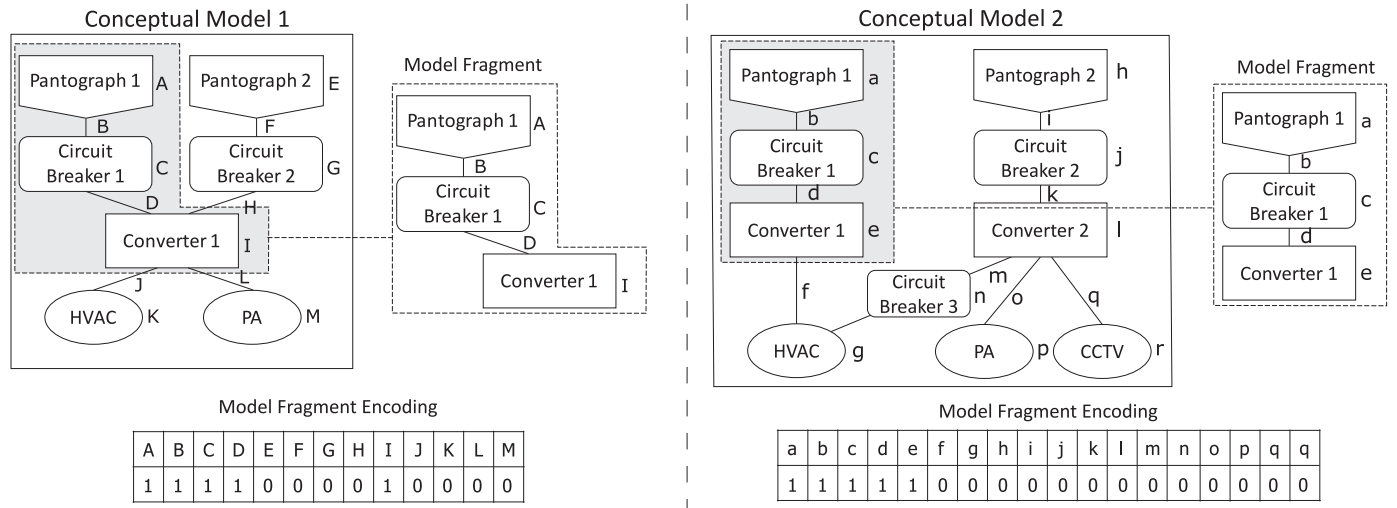


Fig. 3. Examples of Model Fragment Encoding

as its objective. As output, the approach provides a model fragment ranking where each model fragment is ranked taking into account how well the model fragment implements the input requirement.

The top of Fig. 2 shows an example of input to our approach: the model that contains the requirement and the requirement description, which uses natural language to define the target requirement. The center of Fig. 2 shows a simplified representation of the main steps of our approach. Rounded rectangles represent the different steps of the approach, and straight rectangles represent the inputs and outputs of each of the steps. Our approach has three steps:

1. **Initialization:** The first step is to generate a population of model fragments from the model, which serves as input for the evolutionary algorithm. In order to generate the population of model fragments, parts of the model are extracted randomly and added to a collection of model fragments.
2. **Genetic operations:** Second, genetic operations are applied to the model fragment population in order to generate candidate model fragments for the target requirement.
3. **Fitness function:** Finally, the new model fragment population is assessed through the fitness function, which evaluates each of the generated model fragments.

The last two steps of the approach are repeated until the solution converges to a certain stop condition. Usually, the stop condition can be a time slot, a fixed number of generations, or a trigger value of the fitness that makes the process finish when reached (Font et al., 2016b). Since, the stop condition greatly depends on the domain and the problem being solved, it is adjusted depending on the results being output, taking into account when the fitness values are converging and no further improvements are being made by new generations (Font et al., 2016b). When the stop condition is met, the evolutionary algorithm provides a model fragment list, which is ranked according to the objectives for the requirement (see the bottom of Fig. 2).

The following sections describe the genetic operations of TLR-ELtoR for generating new model fragments and how the fitness of each model fragment is determined in terms of similarity to the requirement description.

4. Genetic operations of the TLR-ELtoR approach

The second step of our approach is to generate and to evolve a set of model fragments that could realize the requirement. To do this, this step deals with the encoding of the model fragments and the selection of genetic operators.

4.1. Model fragment encoding for genetic operations

Traditionally, in evolutionary algorithms, each possible solution of the problem is encoded as a string of binary values. However, encoding each model fragment as a string of binary values is not straightforward. The authors in Cetina et al. (2017) propose an encoding where each model fragment is encoded as an individual in relation to the model. In other words, each individual is a set of model elements that are present or absent in a model fragment.

Fig. 3 shows two examples of the representation of model fragments. Each letter labels a model element of the model. Therefore, the individual contains as many positions as model elements in the model and the binary value of these positions depends on the presence or absence of the model elements in the model fragment. If the model element appears in the model fragment, the value will be 1; if the model element does not appear in the model fragment the value will be 0.

Fig. 3 also shows that the encoding will be different for different models, even though the model fragment to be encoded is the same. Both of the examples in Fig. 3 represent the same model fragment. However, since they come from different models, their representations are different.

4.2. Genetic operators

The generation of new model fragments (based on existing ones) is done by applying a set of three genetic operators, which are adapted to work on model fragments. These genetic operations were introduced for the first time in Font et al. (2016b) to carry out the selection of parents, the crossover, and the mutation of model fragments.

The **selection operator** picks the best candidates from the population as input for the rest of the operators. There are different methods that can be used to perform the selection of the parents. One of the most widespread methods (adopted by our work) is to follow the wheel selection mechanism (Affenzeller et al., 2009),

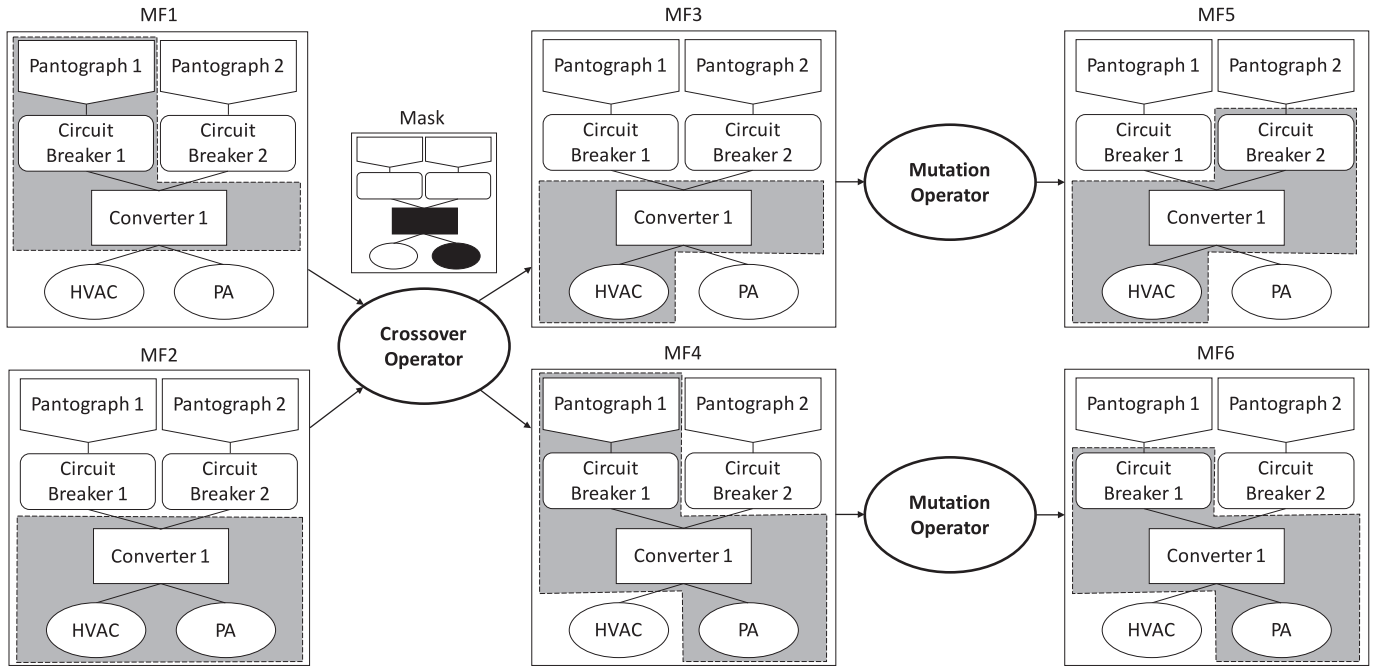


Fig. 4. Example of genetic operations

where each model fragment from the population has a probability of being selected that is proportional to its fitness score. Candidates with high fitness values have higher probabilities of being chosen as parents for the next generation.

The **crossover operation** enables the creation of two new individuals by combining the genetic material from two model fragments. A randomly generated mask determines how the combination is done, indicating for each element of the model fragments if the offspring should inherit from one model fragment or the other. Specifically, the mask is created randomly and all of the model elements have a 50% probability of belonging to the mask. To do this, a random number (0 or 1) is generated for each model element. The elements whose value is 1 belong to the mask and the elements whose value is 0 belong to the inverse of the mask. Moreover, a model fragment is a subset of the elements that are present in a model. Since both model fragments are extracted from the same model, their combination will always return a model fragment that is part of the original product model. As a result of the **crossover operation**, two individuals are generated: one by directly applying the mask, and the other one by applying the inverse of the mask, as is usually done in genetic algorithms [Davis \(1991\)](#).

The **mutation operator** is used to imitate the mutations that occur randomly in nature when new individuals are born. In other words, new individuals have small differences with their parents that could make them adapt better (or worse) to their living environment. Following this idea, the mutation operator applied to model fragments ([Font et al., 2016b](#)) takes as input a model fragment and mutates it into a new one, which is returned as output. Specifically, the mutation operator can perform two kinds of modifications: the addition of elements to the model fragment, or the removal of elements from the model fragment. Since the approach is looking for fragments of the model that realize a specific requirement, the new modified fragment must remain a part of this model. Therefore, the modifications that can be done to the model fragment must be driven by the model, which determines the additions and subtractions of elements that can be applied to the model fragments in the population.

After applying the genetic operators, it may be that not all of the elements of the new individuals are connected. Indeed, the requirement can be implemented by several model elements that are not directly connected in the model ([Font et al., 2016b](#)). Therefore, it is necessary to create fragments of this kind since they could be the ones realizing the requirement.

[Fig. 4](#) shows an example of the application of the two genetic operations. First, the crossover operation is applied. We select the two model fragments to which the operator is applied. Then, the first model fragment (MF1) is combined with the second model fragment (MF2) according to a mask that contains two sets of elements (one regular and one marked in black). To create the first of the new individuals, we interpret the mask by selecting the blackened elements from the first parent (MF1) and the regular elements from the second parent (MF2). As a result, the new model fragment (MF3) contains the set of elements that are present in the mask in MF1 and the set of elements that are absent in the mask in MF2. In addition, the mask is also interpreted in the opposite way by selecting the blackened elements from MF2 and the regular elements from MF1, thus producing another new and distinct model fragment (MF4).

Afterwards, the mutation operator is applied. In this example, the mutation operation takes the first offspring produced through the crossover operator and adds one element (the second circuit breaker). Then, the mutation operation takes the second offspring and removes one element (the first pantograph). The resulting model fragments (MF5 and MF6) are new candidates in the population for the realization of the input requirement.

5. Fitness of the TLR-ELtoR approach

In evolutionary algorithms, this step determines what degree of adaptation to the environment each individual has. Following this idea, in our approach, the fitness step is used to assess how suitable each model fragment is in comparison with the target requirement. To do this, a Learning to Rank algorithm is used to rank a set of model fragments depending on their closeness to a requirement.

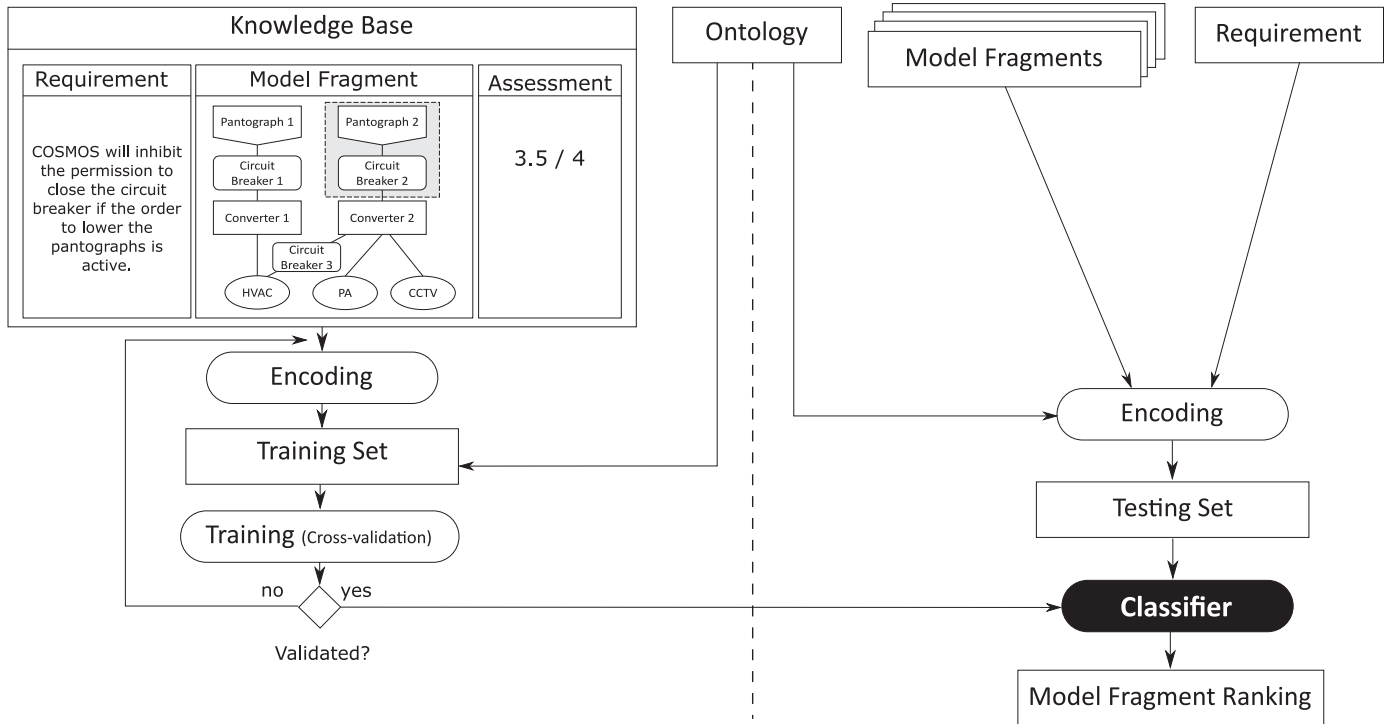


Fig. 5. Overview of the Fitness Function

Learning to Rank (LtoR) is the name given to a family of Machine Learning algorithms, which automatically address ranking tasks. Specifically, the LtoR algorithms make possible the construction of a classifier that contains a set of rules to rank objects. The classifier automatically learns these rules by comparing the objects within a knowledge base. Then, since the classifier knows how to rank objects following these learned rules, the classifier can be used to rank new objects. In other words, Learning to Rank algorithms use a knowledge base to generate a classifier, which is called training. Then, the classifier is used to rank new objects, which is called testing (Lu et al., 2015).

Fig. 5 shows the overview of the Fitness Function, where the LtoR algorithm is applied in our approach. The part on the left of the figure shows the training process where the classifier learns how well each model fragment realizes a specific requirement. To do this, the knowledge base contains traces between requirements and model fragments that are known. The part on the right of the figure shows the testing process where a population of model fragments is ranked by means of the classifier, which determines which model fragment is a better realization of the requirement than another model fragment. Therefore, the classifier is considered as both an artifact (output from the training process) and a step (responsible for ranking in the testing process). For this reason, Fig. 5 shows the classifier in a black, rounded rectangle to point out its double meaning. The following sections provide a more detailed description of the training and testing processes. However, since both training and testing have to encode their inputs (see Fig. 5), the encoding is explained before these processes.

5.1. Model Fragment encoding for the fitness function

Since most of the Machine Learning techniques, such as Learning to Rank algorithms, are designed to process feature vectors as inputs (Bianchini et al., 2013), our model fragments have to be encoded into feature vectors to be able to use LtoR. Feature vectors

are known as the ordered enumeration of features that characterize the object being observed (Chandrashekar and Sahin, 2014).

However, the set of features that are selected to characterize the object have to be the same for all of the model fragments. This guarantees that the feature vectors have the same length and the same features so that the comparison between feature vectors is fair. Moreover, the fitness function determines the suitability between the model fragment and the target requirement. Since both the model fragment and the requirement are being observed, the feature vectors would have to be generated by encoding them both. For these reasons, the encoding for the fitness function is different from the encoding for genetic operations where both the length and the values of the individuals depend on model elements and the requirement is not considered.

In Marcén et al. (2017), we proposed an encoding where each model fragment is encoded as a feature vector taking into account an ontology. Specifically, each concept and relation in the ontology is represented as a feature in the feature vector. The value of each feature is computed as the frequency of the concept or the relation in the model fragment. In the same way, the target requirement is encoded as part of the feature vector taking into account the ontology. Specifically, each concept in the ontology is represented as a feature in the feature vector and the value of each feature is computed as the frequency of the concept in the requirement.

Since both requirements and model fragments are based on natural language, the terms used in the ontology do not always align well with the terms in the requirements and with the terms in the model fragments. For this reason, Natural Language Processing (NLP) techniques are used to process both the requirements and the model fragments before applying the encoding. Specifically, the requirements and the model fragments are processed by a combination of NLP techniques defined in Lapeña et al. (2017), which consists of tokenizing, lowercasing, removal of duplicate keywords, syntactical analysis, lemmatization, and stopword removal.

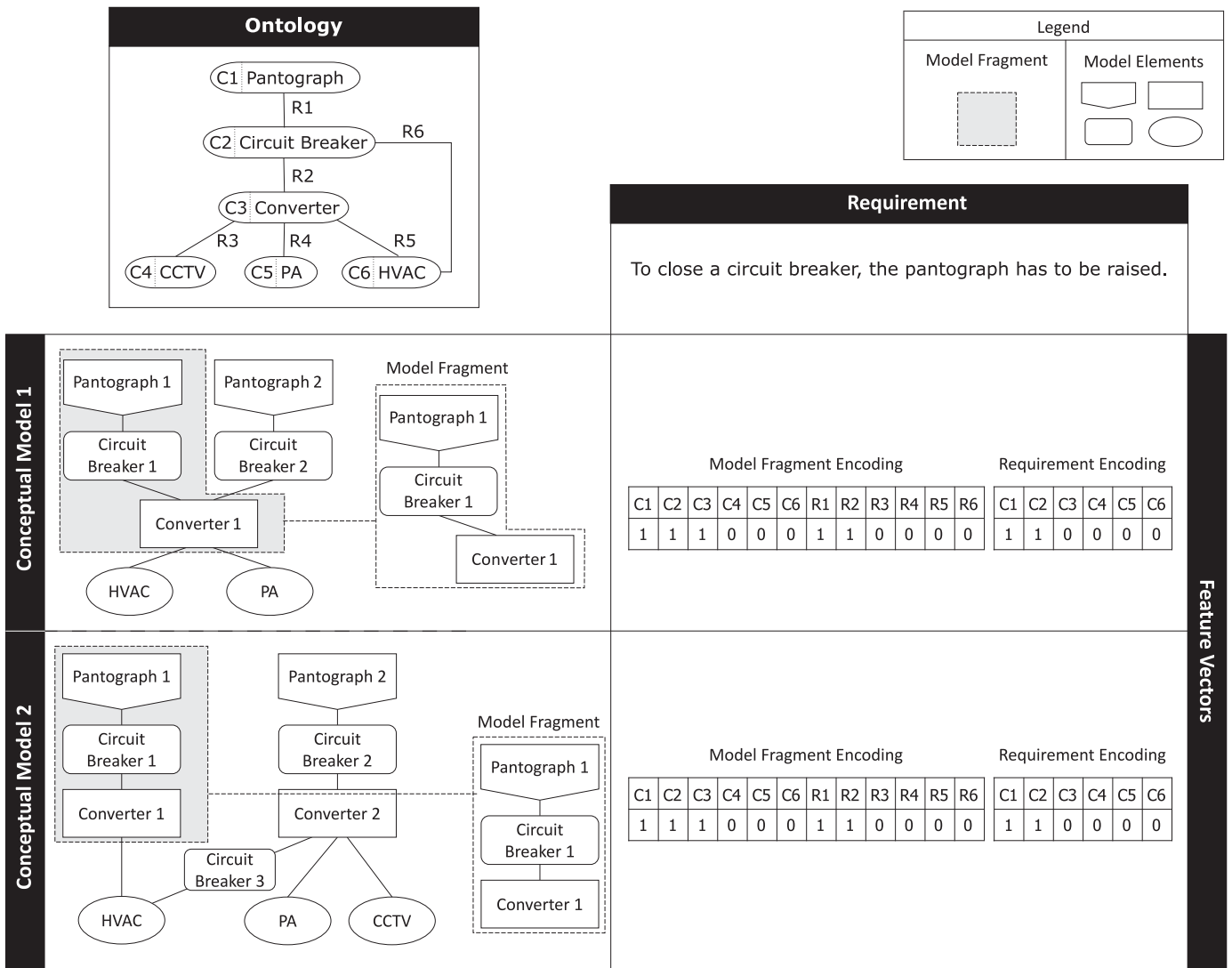


Fig. 6. Examples of Encoding for Model Fragments and Requirements in the Fitness step

Fig. 6 shows two examples of the representation of two model fragments and the target requirement. The concepts and relations of the ontology are features in the feature vector. For example, the concept *Pantograph* is mapped as *C1*, and the relation between the concepts *Converter* and *HVAC* is mapped as *R5*. On the one hand, these concepts and relations are compared with the model fragments, so their values correspond to the number of occurrences of each concept or relation in the model fragment. Therefore, for the first model, the value of the feature *C1* is 1 because there is one pantograph in the model fragment, and the value of the feature *R5* is 0 because there is no relation of the type *Converter-HVAC* in the model fragment. In fact, both models contain a relation of the type *Converter-HVAC*, but none of the model fragments contain this relation. Therefore, the value for *R5* is 0 in both feature vectors because each feature vector only contains the encoding of a model fragment, not the encoding of the whole model. On the other hand, the concepts are also compared with the target requirement, so their values correspond to the number of occurrences of each concept in the requirement. Therefore, the value of the *C1* is 1 because the concept *pantograph* appears once in the requirement. Fig. 6 also shows that the feature vectors do not depend on whole models because they only represent the encoding of model fragments and requirements. Therefore, two model fragments may be

result in the same feature vector although their models are different.

5.2. Training process

The target of the training process is to produce a classifier from a training set, which ranks the model fragments generated by the genetic operations. To do this, the knowledge base has to be encoded to obtain the training set, which is used in a LtoR algorithm to generate a classifier. Also, before using this classifier in the testing process, the validation of its performance is a good practice in order to improve the results.

The knowledge base is composed of traces between requirements and model fragments that are known. Specifically, the knowledge base consists of a set of traces that are generated using the domain experts' experience, results, and documentation, where each trace of the knowledge base is composed of a requirement, a model fragment, and an assessment. The requirement uses natural language to define the requirement. The model fragment consists of an element or a set of elements that belongs to a model. The assessment determines if the model fragment realizes the requirement to a greater or lesser extent. In other words, the assessment determines the similarity between the requirement and the model

fragment. Fig. 5 shows an example of the knowledge base for performing requirement traceability.

In order to apply LtoR algorithms in models, the first step consists of encoding the traces of the knowledge base into the feature vectors. Therefore, the model fragment and the requirement of each trace are encoded following the encoding for the fitness function (see Section 5.1), and the assessment is also included as part of the encoding. Then, the obtained feature vectors compose the training set.

The training set is used to train a classifier, which learns a rule-set through the comparison of the feature vectors of the training set (Shabtai et al., 2009). However, before using this classifier to rank the model fragment in the testing process, it is worth analyzing the performance of the classifier through cross-validation. Cross-validation is a statistical method of evaluating and comparing ML algorithms by dividing data into two segments: one used to train a classifier, and the other used to validate the classifier (Refaeilzadeh et al., 2009). Moreover, to reduce variability, multiple rounds of cross-validation are performed using different partitions, and the results are averaged over the rounds (Song et al., 2011).

The results of the cross-validation provide the performance of the classifier. If this performance is not considered suitable, it is necessary to perform another training iteration. In this iteration, some artifacts of the training process (e.g., the encoding, the ontology, the knowledge base, or the LtoR algorithm) have to be modified in order to improve the classifier. Otherwise, if the performance is considered suitable, the classifier obtains the go-ahead, so the classifier trained with the whole knowledge base is used in the testing process. Once the classifier has been generated, the training process does not have to be repeated again. The same classifier is used whenever the Fitness Function is applied. Therefore, the training process is only performed in the first iteration of the EA, when there is not yet a classifier.

5.3. Testing process

In our approach, the classifier is used to rank the model fragments that are generated after each iteration of the EA. Specifically, the classifier assigns a score to each model fragment based on its closeness to the requirement. For example, taking into account the requirement of Fig. 6, the model fragments must contain at least one *circuit breaker* so that these model fragments are close to this requirement. Therefore, the model fragment with a *circuit breaker* would obtain a better fitness score than a model fragment without a *circuit breaker*.

However, before ranking the model fragments, both the model fragments and the requirement have to be encoded into feature vectors (see Section 5.1) so that the classifier can understand them. The feature vectors obtained from the encoding compose the testing set. Then, each feature vector of the testing set is tested by the classifier, which used the learned rule-set in the training process to assign a fitness score to each one of them.

The fitness score is a numerical value that is greater than 0. If the fitness score is close to 0, the model fragment is not close to the requirement, so the model fragment is not relevant to the requirement. In contrast, the greater the fitness score, the more relevant the model fragment is. Taking into account the fitness scores, the model fragments can be ordered in a ranking where the top positions are occupied by the model fragments with the highest relevance to the requirement.

Finally, as a result, this ranking of model fragments is returned by the Fitness Function, (see Fig. 2). Therefore, in each iteration of the EA, the Fitness Function provides a ranking of model fragments organized by their fitness score. Then, if the stop condition is satisfied, this ranking of model fragment is obtained as a result of the approach. However, if the stop condition is not satisfied, a new

iteration is performed. Therefore, the model fragments are considered as the new population to be evolved; then the genetic operations select the best model fragments, mutate them, and create new model fragments from them.

6. Evaluation

This section presents the evaluation of our approach: the experimental setup, the baselines, a description of the case study where we applied the evaluation, the implementation details, and the obtained results.

6.1. Experimental setup

The goal of this experiment is to perform TLR between requirements and models through TLR-ELtoR and to compare the results with the TLR approaches that have obtained the best results in the literature. Fig. 7 shows an overview of the process that was followed to evaluate our approach (TLR-ELtoR) and the baselines (TLR-Linguistic, TLR-IR, TLR-FNN, TLR-RNN, TLR-LtoR). The top part of Fig. 7 shows the inputs, which are extracted from the documentation provided by our industrial partner: knowledge base, ontology, requirements, product models, and approved traceability between requirements and product models. Each test case is comprised of a requirement, a model of a product, the ontology, and the knowledge base. However, the ontology and the knowledge base are ignored by TLR-Linguistic and TLR-IR because they do not need it. The Oracle is composed of the approved traceability between the requirements and the models.

For each test case, our approach generates a ranking of model fragments. Each model fragment contains the elements of the model that are related to the requirement, so each model fragment fits the traceability between the model and the requirement to a greater or lesser extent. Then, we take the best solution of the ranking and compare it against the oracle, which is the ground truth. Once the comparison is performed, a confusion matrix is calculated. The baselines also recover the traceability links between the requirement and the model in the test cases. As a result, each baseline generates a model fragment. These model fragments are also compared with the oracle, and a confusion matrix is calculated for each baseline. Therefore, we obtain six confusion matrices, one for our TLR-ELtoR approach and one for each baseline.

A confusion matrix is a table that is often used to describe the performance of a classification model (in this case, both the TLR-ELtoR and the baselines) on a set of test data (the solutions) for which the true values are known (from the oracle). In our case, each solution that is outputted by the approaches is a model fragment that is composed of a subset of the model elements that are part of the product model. Since the granularity is at the level of model elements, the presence or absence of each model element is considered as a classification. The confusion matrix distinguishes between the predicted values and the real values, classifying them into four categories:

- True Positive (TP): values that are predicted as true (in the solution) and are true in the real scenario (the oracle).
- False Positive (FP): values that are predicted as true (in the solution) but are false in the real scenario (the oracle).
- True Negative (TN): values that are predicted as false (in the solution) and are false in the real scenario (the oracle).
- False Negative (FN): values that are predicted as false (in the solution) but are true in the real scenario (the oracle).

Then, some performance measurements are derived from the values in the confusion matrix. Specifically, we create a report that includes four performance measurements (recall, precision, the F-

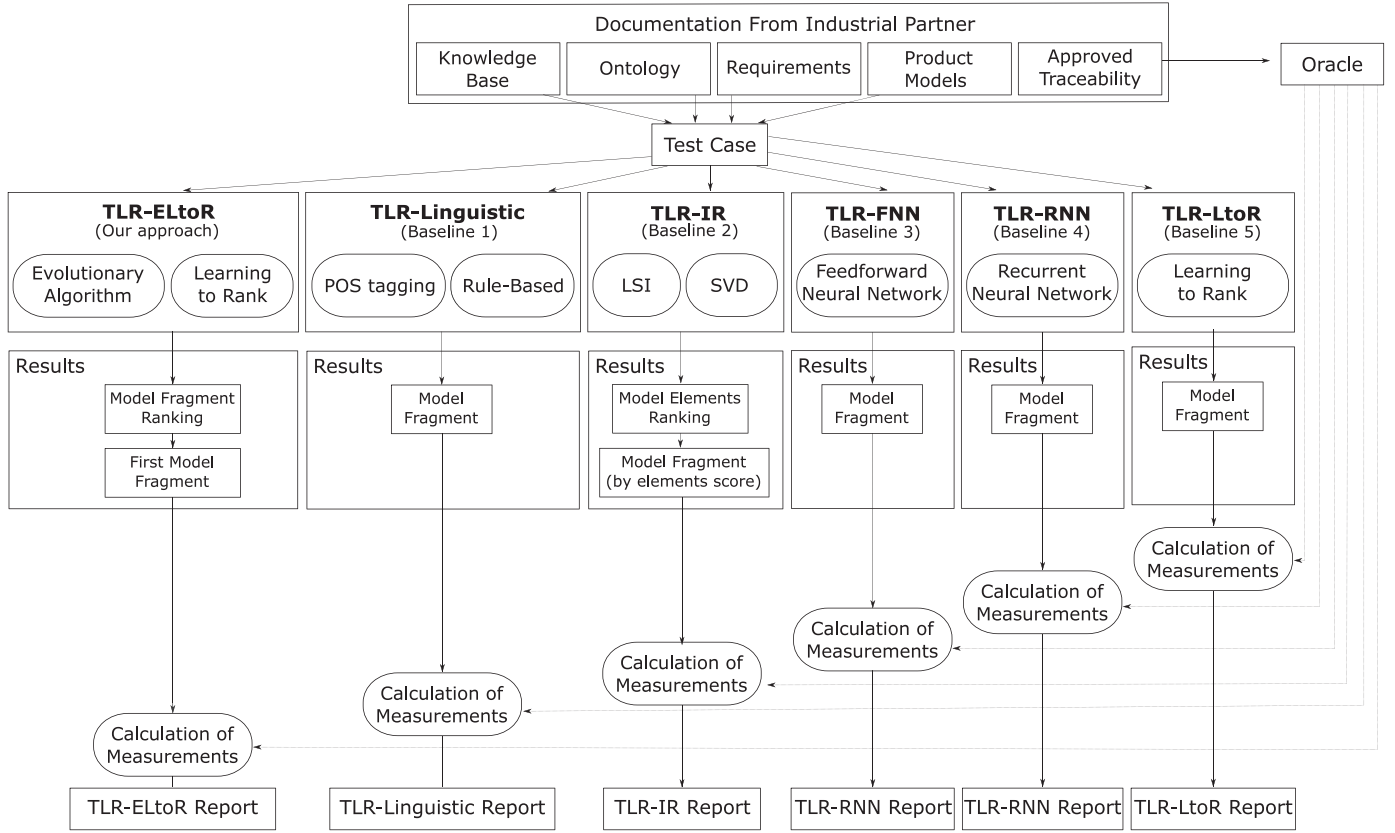


Fig. 7. Experimental setup

measure, and the Matthews Correlation Coefficient) for the test case for both the TLR-ELtoR and the baselines.

Recall measures the proportion of elements of the solution that are correctly retrieved by the proposed solution and is defined as follows:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Precision measures the proportion of elements from the solution that are correct according to the ground truth (the oracle) and is defined as follows:

$$\text{Precision} = \frac{TP}{TP + FP}$$

The F-measure corresponds to the harmonic mean of precision and recall and is defined as follows:

$$F - \text{measure} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 * TP}{2TP + FP + FN}$$

However, none of these previous measures correctly handle negative examples (TN). The **MCC** is a correlation coefficient between the observed and predicted binary classifications that takes into account all of the observed values (TP, TN, FP, FN) and is defined as follows:

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Recall values can range between 0% (i.e., no single model element from the realization of the requirement obtained from the oracle is present in the model fragment of the solution) and 100% (i.e., all of the model elements from the oracle are present in the solution). Precision values can range between 0% (i.e., no single model element from the solution is the oracle) and 100% (i.e., all of the model elements from the solution are present in the oracle). A

value of 100% precision and 100% recall implies that both the solution and the requirement realization from the oracle are the same. MCC values can range between -1 (i.e., there is no correlation between the prediction and the solution) and 1 (i.e., the prediction is perfect). Moreover, a MCC value of 0 corresponds to a random prediction.

6.1.1. TLR-ELtoR setup

This section describes the technical details of our TLR-ELtoR approach taking into account the experimental setup defined. Specifically, four technical details are addressed: the stop condition, the hyperparameters for the evolutionary algorithm, the LtoR algorithm with its setting parameters, and the cross-validation method.

In general, there are two atomic performance measures for evolutionary algorithms: one regarding solution quality, and one regarding algorithm speed or search effort. In this paper, we focus on the solution quality (i.e., obtaining a solution that is more similar to the one from the oracle in terms of precision and recall). After running some prior tests to determine the number of iterations to converge (and adding a margin to ensure convergence), we allocated a fixed amount of iterations (200 iterations) to stop the execution.

For the settings of the evolutionary algorithm, namely population size, crossover probability, and mutation probability, we have chosen the values 100, 0.9, and 0.1, respectively. These were selected based on the parameters that are commonly used in the literature (Sayyad et al., 2013) and the results of some preliminary tuning experiments.

With regard to the LtoR algorithm, the selection of this algorithm depends on several aspects, such as the size of the knowledge base. RankBoost (Freund et al., 2003) belongs to the family of LtoR and is well known for its efficiency and effectiveness in different domains (Canuto et al., 2013; Cao et al., 2018). Moreover,

RTOM_RULE Rule-1: EXISTS SEQUENCE(<x1/{NN1, NN2}>, <x2/{VBZ, VBR}>, <x3/{JJ}>) in Requirement; <x4/CLASS>, <x5/ATTRIBUTE> in Model SUCH THAT ATTRIBUTE_OF(<x5>, <x4>) and CONTAINS(NAME(<x5>), <x3>) and (CONTAINS(NAME(<x4>), <x1>) or CONTAINS(NAME(<x4>), SINGULAR_FORM<x1>)) ACTION GENERATE OVERLAPS(Requirement, <x5>) RTOM_RULE_END
RTOM_RULE Rule-2: EXISTS SEQUENCE(<x1/{VV0, VVI, VVZ}>, <x2/{AT}>, <x3/{NN1, NN2}>) in Requirement; <x4/CLASS>, <x5/OPERATION> in Model SUCH THAT OPERATION_OF(<x5>, <x4>) and MEMBER_OF(<x1>, SYNONYMS(STEREOType(<x5>))) and CONTAINS(NAME(<x4>), <x3>) or CONTAINS(NAME(<x4>), SINGULAR_FORM<x3>)) ACTION GENERATE REQUIRES_EXECUTION_OF(Requirement, <x4>) RTOM_RULE_END

Fig. 8. Example of requirement-to-object-model rules

Rankboost can benefit from a small knowledge base together with a small number of features in the encoding to reduce the overfitting problem (Zhou and Feng, 2017; Wang et al., 2014). Since this condition is satisfied by our case study, TLR-ELtoR was guided by Rankboost with the parameters tuned as in Kırac et al. (2018). First, a grid search was built to determine the values of the parameters: number of iterations, and metric. Then, we uniformly sampled each of the two parameters in their range and evaluated all of the combinations of the sampled values. As a result, the parameters were tuned with *iteration* = 200 and *metric* equal to ERR10.

Moreover, even though our approach considers cross-validation as a step of the Fitness Function, the approach does not restrict the possibilities to one specific cross-validation method. In this evaluation, TLR-ELtoR used a *k*-fold validation with a *k* value equal to 4. The *k*-fold validation is the most popular cross-validation procedure. Specifically, this method consists of randomly dividing the knowledge base into *k*-independent partitions. Then, *k* – 1 of the partitions are used to train the classifier, and this classifier is then used to test the partition that is left out. This procedure is repeated *k* times, each time leaving out another partition. This produces *k* estimations of the classifier, allowing assessment of its central tendency and variance (Hirzel et al., 2006).

6.2. Baselines

Winkler and Pilgrim (2010) classify several approaches that have been created over the past 15 years that try to optimize the automatic identification of traces. Based on this classification, as baselines, we selected the two approaches that obtain the best results for traceability links between requirements and models: (1) a rule-based approach that deduces traces by applying rules (TLR-Linguistic) (Spanoudakis et al., 2004); and (2) an information retrieval approach that can detect candidate traceability links through Information Retrieval (TLR-IR) (De Lucia et al., 2004; Lucia et al., 2007).

Deep learning techniques have also successfully been applied in TLR in some recent works (Guo et al., 2017). Therefore, we decided to compare our approach with two baselines that apply deep learning: (1) the first one is based on a Feedforward Neural Network (TLR-FNN); and (2) the second one is based on a Recurrent Neural Network (TLR-RNN).

Finally, to check the need for the evolutionary algorithm in our approach, TLR-ELtoR is also compared to TLR-LtoR, which explores

the search space by means of brute-force. Therefore, the model fragments are generated from the model and evaluated through LtoR, but the results obtained from the LtoR process are not used to guide the generation of new model fragments. Since there is no guide to explore the model, the search for the model fragment that realizes a specific requirement is performed by brute-force.

6.2.1. TLR-Linguistic: linguistic rule-based baseline

Spanoudakis et al. (2004) present a linguistic rule-based approach to support the automatic generation of traceability links between requirements and models. Specifically, the traceability links are generated following two stages:

- Stage 1: a Parts-of-Speech (POS) tagging technique (Leech et al., 1994) is applied on the requirements that are defined using natural language.
- Stage 2: the traceability links between the requirements and the models are generated through the *requirement-to-object-model* rules.

The *requirement-to-object-model* (RTOM) rules are specified by investigating grammatical patterns in requirements. Moreover, the RTOM rules are based on two kinds of relations between requirements and models. On one hand, *Overlap* relations are understood to be the relation between a sequence of terms in a requirement and a class, attribute, association, or association-end in a model. On the other hand, *Requires_Execution_Of* relations are understood to be the relation between a sequence of terms in a requirement and an operation in a model.

Fig. 8 shows an example of both kinds of rules following the syntax that is defined in Spanoudakis et al. (2004). The top rule can establish an *Overlap* relation between a requirement and an attribute in a model. The bottom rule can establish a *Requires_Execution_Of* relation between a requirement and an operation in a model. These rules generate the traceability links between the requirement and the model presented in Fig. 9.

The first rule in Fig. 8 attempts to match a syntactic expression that consists of a noun (< x1/{NN1, NN2} >), the verb *to be* in the present form (< x2/{VBZ, VBR} >), and an adjective (< x3/{JJ} >) with an attribute in the model. The matching succeeds if: (a) the name of the attribute contains the adjective and the name of the class that defines the attribute contains the noun; or (b) the name of the attribute contains the adjective and

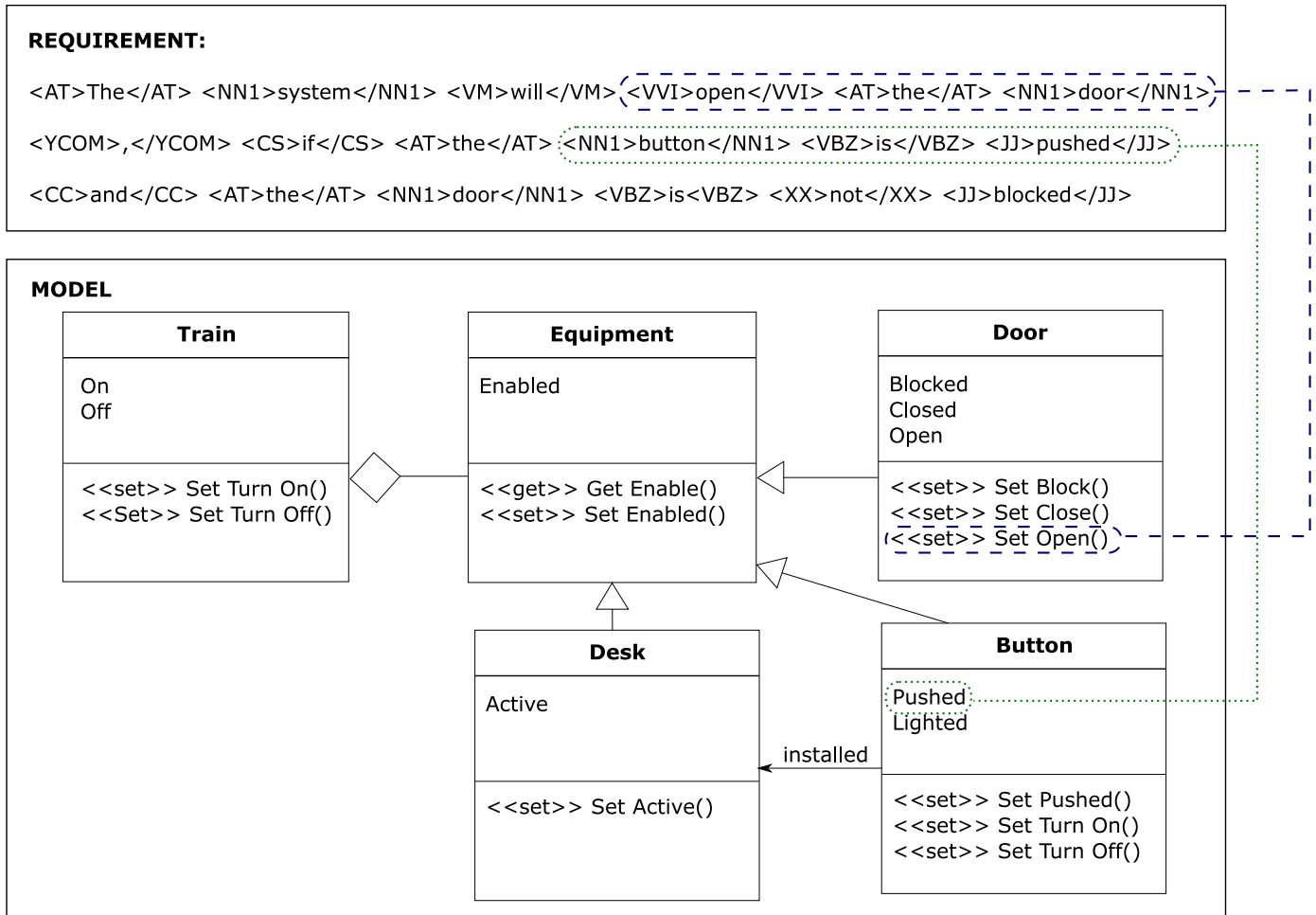


Fig. 9. Example of traceability links generation based on RTOM rules

the name of the class that defines the attribute contains the singular form of the noun. Therefore, in Fig. 9, the sequence of terms $\langle \text{NN1} \rangle \text{button} \langle / \text{NN1} \rangle \langle \text{VBZ} \rangle \text{is} \langle / \text{VBZ} \rangle \langle \text{JJ} \rangle \text{pushed} \langle / \text{JJ} \rangle$ in the requirement and the attribute Pushed of the class Button satisfy the conditions of the rule. As a consequence, an *Overlap* relation is created between them.

The second rule in Fig. 8 attempts to match a syntactic expression that consists of a verb ($\langle x1/\{\text{VVO}, \text{VVI}, \text{VVZ}\} \rangle$), an article ($\langle x2/\{\text{AT}\} \rangle$), and a noun ($\langle x3/\{\text{NN1}, \text{NN2}\} \rangle$) with an operation in the model. The matching succeeds if: (a) the name of the operation contains the verb or is a synonym of the verb and the name of the class of the operation contains the noun; or (b) the name of the operation contains the verb or is a synonym of the verb and the name of the class of the operation contains the singular form of the noun. Therefore, in Fig. 9, the sequence of terms $\langle \text{VVI} \rangle \text{open} \langle / \text{VVI} \rangle \langle \text{AT} \rangle \text{the} \langle / \text{AT} \rangle \langle \text{NN1} \rangle \text{door} \langle / \text{NN1} \rangle$ in the requirement and the operation Set Open of the class Door satisfy the conditions of the rule. As a consequence, a *Requires_Execution_Of* relation is created between them.

In Spanoudakis et al. (2004), there are two different types of traceability rules: RTOM for traceability relations between requirements and model elements, and inter-requirement rules for traceability relations between different parts of a requirement statement. In total, the authors propose 26 rules for two domains: a software-intensive TV system created by Philips, and a university course management system. Since our approach is focused only on

the traceability between requirements and model elements, this baseline only tackles the RTOM traceability rules for our domain. Therefore, based on the guides and the examples of rules that are provided by Spanoudakis et al. (2004), a domain expert who was not involved in the research generated an initial set of rules for our domain. In addition, to mitigate the dependence on a single domain expert, a second expert who also was not involved in the research extended the set of rules. In the end, the extended set contains nine RTOM rules, which is similar to the number proposed by Spanoudakis et al. (2004). However, there is no significant difference between the results obtained using the initial set and the results obtained using the extended set. Specifically, the results described in this work correspond to the extended set, which are a bit better than those obtained from the initial set. Nonetheless, in both cases, the results are not as good as the ones obtained with our approach.

6.2.2. TLR-IR: Information Retrieval baseline

Information Retrieval (IR) (Frakes and Baeza-Yates, 1992; Manning et al., 2008; Salton and McGill, 1986) is a sub-field of computer science that deals with the automated storage and retrieval of documents. IR techniques have been successfully used to retrieve traceability links between different kinds of software artifacts in different contexts (Lucia et al., 2000; Oliveto et al., 2010; Antoniol et al., 2002; Marcus and Maletic, 2003a; De Lucia et al., 2006). Specifically, in De Lucia et al. (2004) and Lucia et al. (2007), De Lucia et al. use Latent Semantic Indexing (LSI) to recover traceability links between requirements and different kinds of soft-

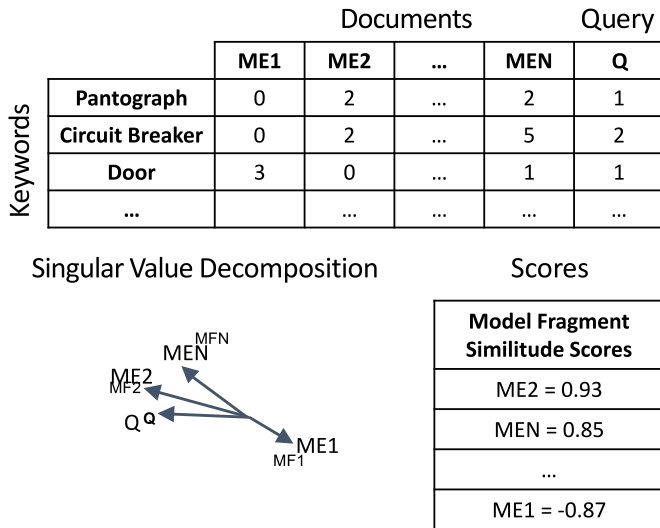


Fig. 10. Example of traceability link recovery using latent semantic indexing

ware artifacts, including models in the form of use-case diagrams, among others. We use LSI to recover traceability links between requirements and models as one of the baselines for our work.

Specifically, given a certain requirement-model pair as input for LSI, we use the produced outcome of the technique to build a model fragment from the model that serves as a candidate for realizing the requirement. The following paragraphs provide more details on the process.

Latent Semantic Indexing (LSI) (Landauer et al., 1998) is an automatic mathematical/statistical technique that analyzes relationships between queries and documents (bodies of text). Since both queries and documents are based on natural language, Natural Language Processing (NLP) techniques are used to process them. In fact, NLP has a direct and beneficial impact on the results, so before applying LSI, the queries and the documents are processed by a combination of NLP techniques defined in Lapeña et al. (2017), which consists of tokenizing, lowercasing, removal of duplicate keywords, syntactical analysis, lemmatization, and stopword removal. Then, LSI constructs vector representations of both a user query and a corpus of text documents by encoding them as a term-by-document co-occurrence matrix and analyzes the relationships between those vectors to get a similarity ranking between the query and the documents (see Fig. 10).

Fig. 10 shows an example term-by-document co-occurrence matrix, with values associated to our case study, the vectors, and the resulting ranking. An overview of the elements of the matrix is provided in the following paragraphs:

- Each row in the matrix (*term*) stands for each of the words that compose the processed requirement and NL representation of the input model. The NL representation of the input model is extracted using the technique presented in Meziane et al. (2008). For example, Fig. 10 shows a set of representative words in the domain such as 'pantograph' or 'door' as the *terms* of each row.
- Each column in the matrix (*document*) stands for one model element from one input model, taken from our real-world case study. For example, Fig. 10 shows identifiers in the columns such as 'ME1' or 'ME2', which stand for the *documents* of those specific model elements.
- The final column stands for the *query*, which is a requirement in our case study.
- Each cell in the matrix contains the frequency with which the *term* of its row appears in the *document* denoted by its column.

For instance, in Fig. 10, the term 'pantograph' appears twice in the 'ME2' document and once in the query.

Vector representations of the documents and the query are obtained by normalizing and decomposing the term-by-document co-occurrence matrix using a matrix factorization technique called Singular Value Decomposition (SVD) (Landauer et al., 1998). SVD is a form of factor analysis, or more properly, it is the mathematical generalization of which factor analysis is a special case. In SVD, a rectangular matrix is decomposed into the product of three other matrices. One component matrix describes the original row entities as vectors of derived orthogonal factor values, a second one describes the original column entities in the same way, and the third one is a diagonal matrix containing scaling values such that when the three components are matrix-multiplied, the original matrix is reconstructed.

A three-dimensional graph of the SVD is provided in Fig. 10. The graph shows the vectorial representations of some of the matrix columns. For legibility reasons, only a small set of the columns is represented. To measure the degree of similarity between vectors, the cosine between the query vector and the documents vectors is calculated. Cosine values that are closer to 1 denote a higher degree of similarity, and cosine values that are closer to -1 denote a lower degree of similarity. Similarity increases as vectors point in the same general direction (as more terms are shared between documents). With this measurement, the model elements are ordered according to their degree of similarity to the requirement.

The relevancy ranking (which can be seen in Fig. 10) is produced according to the degrees of similarity calculated. In this example, LSI retrieves 'ME2' and 'MEN' in the first and second position of the relevancy ranking since the query-documents cosines are '0.9343' and '0.8524', implying a high degree of similarity between the model elements and the requirement. In contrast, the 'M1' model element is returned to a lower position in the ranking since its query-document cosine is '-0.8736', implying a lower degree of similarity.

From the ranking, of all the model elements, only those model elements that have a degree of similarity greater than x must be taken into account. A good heuristic that is widely used is $x = 0.7$. This value corresponds to a 45° angle between the corresponding vectors. Even though the selection of the threshold is an issue under study, the heuristic chosen for this work has yielded good results in other similar works (Marcus et al., 2004; Salman et al., 2014).

Following this principle, the elements with a degree of similarity equal or greater than to $x = 0.7$ are taken to conform a model fragment, which is a candidate for realizing the requirement. In the example provided in Fig. 10, ME2 and MEN are model elements that conform part of the model fragment that is obtained by this baseline for the requirement because their cosine values are greater than the 0.7 threshold. The model fragment generated in this manner is the final output of the TLR-IR baseline.

6.2.3. TLR-FNN: Feedforward Neural Network baseline

Feedforward Neural Networks (FNNs) represent a traditional neural network structure and lay the foundation for many other structures (Haykin, 1994). Data flow always moves one direction, from input layer to hidden layer, then to output layer; it never goes backwards. Fig. 11 shows the structure of a FNN where the FNN receives a vector of I input signals, $z = (z_1, z_2, \dots, z_I)$. The neurons of the hidden layer assign to each input signal, z_i , its respective weight, v_i , to strengthen or deplete the input signal. Weighted inputs are accumulated at each neuron and then an activation function determines the output (or firing strength) of each neuron, o . In fact, the strength of the output is further influenced by a threshold value, which is also referred to as the bias; thus, the activation

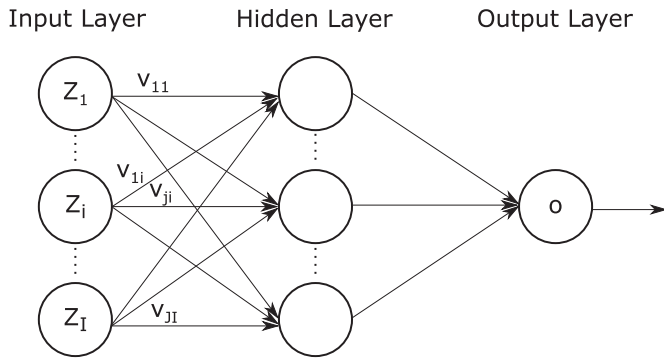


Fig. 11. Feedforward neural network

function receives both the input signal and the bias to determine the output of each neuron (Engelbrecht, 2007).

While Fig. 11 shows only one hidden layer, a FNN can have more than one hidden layer. However, it has been proved that FNNs with monotonically increasing differentiable functions can approximate any continuous function with one layer, provided that the hidden layer has enough hidden neurons (Hornik et al., 1989). Specifically, the network architecture of the FNN implemented here is a dense layer that is followed by the final softmax layer. Moreover, we performed a hyperparameter optimization based on the random search optimization provided by the Deep Learning for Java library. For all of the layers, the hyperparameter optimization resulted in an initial learning rate of 0.0035, and the Gaussian dis-

tribution recommended in Klambauer et al. (2017) for weight initialization. In addition, for the dense layer, the hyperparameter optimization resulted in a layer size of 128 and the randomized rectified linear unit (RReLU) as the activation function.

In addition, since FNN is a typical kind of supervised machine learning method, the training process is required to adjust weights and bias for inputs. Therefore, a training set has to be used to train the FNN and to determine how well the FNN has learned (Zhang et al., 0000). Moreover, the testing process is required to recover the traceability link between a model and a requirement. However, both the training set and the testing are a bit different from the sets used in TLR-ELtoR. Since this baseline is not based on a EA, it cannot generate and evolve model fragments. Therefore, this baseline works at the model-element level instead of at the model-fragment level. For this reason, both the training set and the testing set consist of a set of vectors, where each vector represents the relation between a model element and a requirement.

Fig. 12 shows an example of the encoding of the knowledge base at the model-element level. From a sample of the knowledge base, each element of the model is encoded by means of a feature vector. Then, the encoding for the requirement is included in all of the feature vectors. Finally, a numerical value is included as target for the training. This value is equal to 0, if the model element is not present in the model fragment. Otherwise, this value is equal to 1, if the model element is present in the model fragment. This value is used by the FNN to learn relations between the model elements and the requirements. In summary, each sample of the knowledge base is encoded using several feature vectors, one for each element in the model.

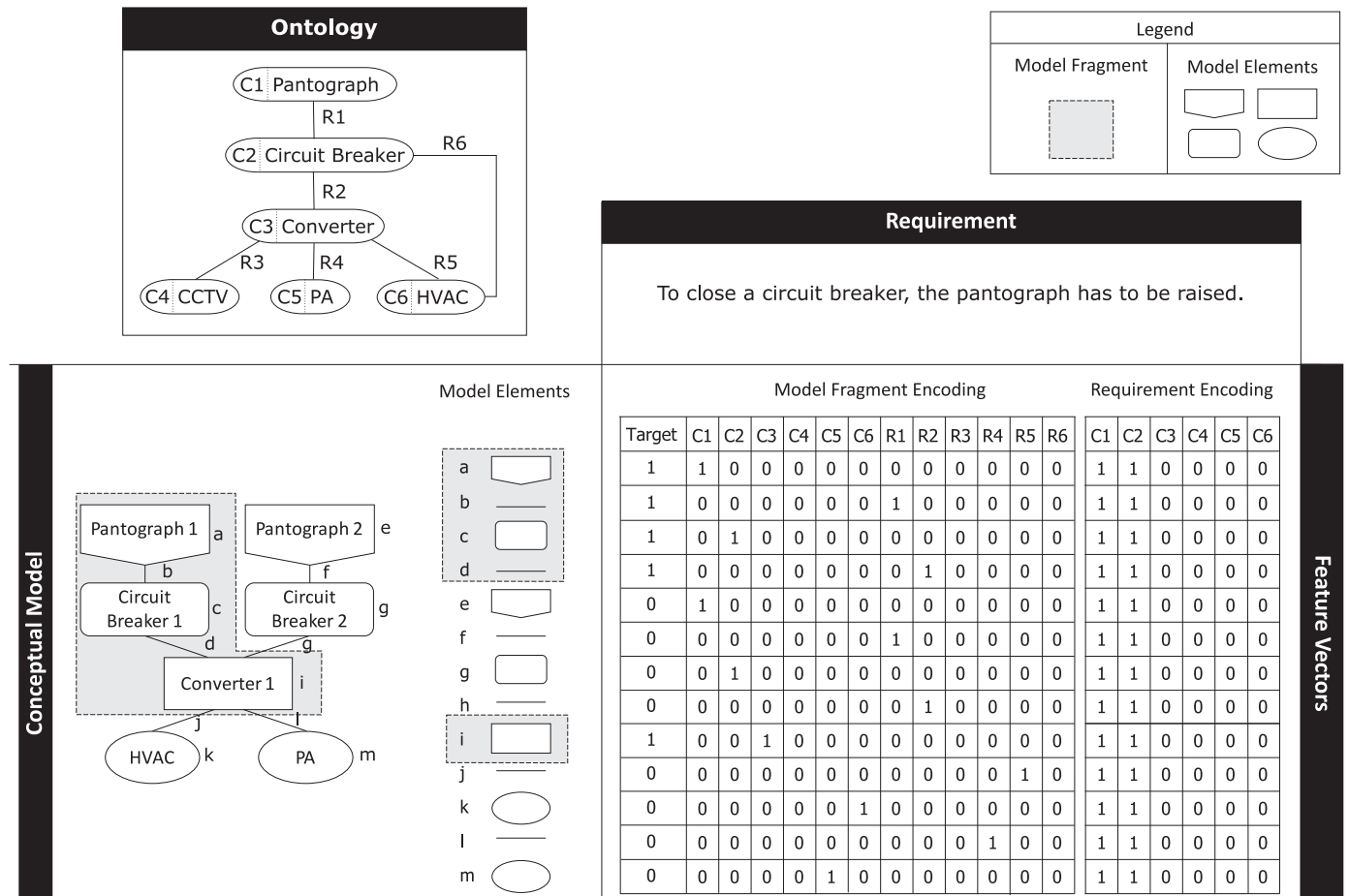


Fig. 12. Example of encoding of the knowledge base at the model-level

Likewise, the model and the requirement for the testing set are encoded at the model-element level. Each model element leads to one feature vector, and the encoding of the requirement is included in all of the feature vectors.

In addition, Fig. 12 also shows a limitation of the encoding proposed in Marcén et al. (2017) when it is applied at the model-element level. Several feature vectors contain the same feature values, but different target values (e.g., both pantographs). Taking into account the ontology in Fig. 12, there is no way to differentiate between two elements of the same type; for example, both pantographs have the same encoding. However, the ontology can be extended to tackle specific properties of each type of element. For example, the ontology of Fig. 12 can be extended to include the status of each pantograph as an attribute of the pantograph concept. Therefore, the two pantographs can be differentiated taking into account if their status is in the up, down, or middle position. Specifically, in our case study, the ontology was extended with 14 attributes that empowered us to mitigate this threat.

6.2.4. TLR-RNN: Recurrent Neural Network baseline

Since the number of parameters in a fully connected FNN can grow extremely large as the width and depth of the network increases, researchers have proposed other neural network structures targeting different types of practical problems. Recurrent Neural Networks (RNNs) are particularly well suited for processing sequential data such as text and audio. While FNNs have no feedback connections to previous layers, RNNs have these feedback connections to model the temporal characteristics of the problem being learned (Engelbrecht, 2007). Moreover, RNNs have successfully been applied in TLR in some recent works (Guo et al., 2017).

Although RNNs are specifically designed to process sequential data, RNNs have showed great results in some cases of non-sequential input information, for instance, image captioning (Mao et al., 2014) or prediction of hospital readmission (Chopra et al., 2017). In these works, even if the input data is not in the form of sequences, they can make classifiers able to learn so that they process data in sequential order only (Chopra et al., 2017). In our case, even if the models are not sequential data, we can order the feature vectors of the model elements so that a classifier trained by a RNN benefits from the sequential order of the model elements. For example, taking into account the Fig. 12, the three first feature vectors match the model elements: *Pantograph1*, the relation *Pantograph1-Circuit Breaker1*, and *Circuit Breaker1*. Taking into account their order, even if we knew that *Pantograph1* is related to the requirement, we could not determine if the other two model elements are related to the requirement. However, if we knew that the *Pantograph1* and the relation *Pantograph1-Circuit Breaker1* are related to the requirement, it would be certainly reasonable to assume that the *Circuit Breaker1* is related to the requirement.

Fig. 13 shows the structure of an Elman RNN, which is a RNN based on the extension of a FNN. As illustrated in Fig. 13, data flow moves from an input layer to a hidden layer, but there is a new layer, named context layer, that makes a copy of the hidden layer. This context layer serves as an extension of the input layer, feeding signals that represent previous network states to the hidden layer. Therefore, the input vector is $z = (z_1, \dots, z_l, z_{l+1}, \dots, z_{l+j})$, where the first l signals are the actual inputs of the network and the j signals are the context units (Engelbrecht, 2007).

A prominent drawback of the standard RNN model is that the network degrades when long dependencies exist in the sequence due to the phenomenon of exploding or vanishing gradients during back-propagation (Bengio et al., 1994). This makes a standard RNN model difficult to train. The exploding gradient problem can be effectively addressed by scaling down the gradient when its norm is bigger than a preset value (i.e., Gradient Clipping) (Bengio et al.,

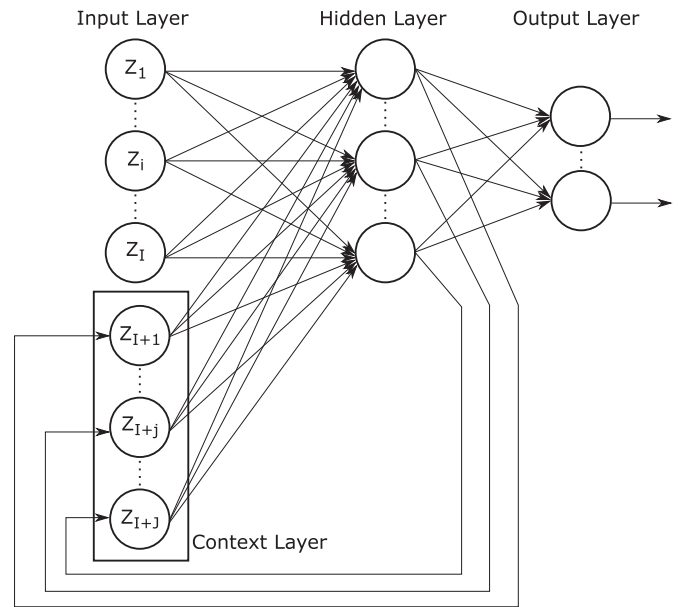


Fig. 13. Elman Simple Recurrent Neural Network

1994). To address the vanishing gradient problem of the standard RNN model, the RNN network that is used as baseline in this work applies Long Short Term Memory (LSTM), which is a variant provided by researchers that has mechanics to preserve long-term dependencies (Guo et al., 2017).

LSTM networks include a memory cell vector in the recurrent neuron to preserve long-term dependencies (Hochreiter and Schmidhuber, 1997). LSTM also introduces a gating mechanism to control when and how to read or write information to the memory cell. A gate in LSTM usually uses a sigmoid function $\sigma(z) = 1/(1 + e^{-z})$ and controls information throughput using a point-wise multiplication operation. Specifically, when the sigmoid function outputs 0, the gate forbids any information from passing, while all information is allowed to pass when the sigmoid function output is 1 (Guo et al., 2017). Each LSTM neuron contains an input gate, a forget gate, and an output gate. The input gate controls how much each signal in a candidate vector should be "remembered". The forget gate controls how much each signal in the previous memory neuron state should be retained, so the neuron "remembers" information until it is erased by the forget gate. Finally, the output gate controls when a signal output is used in the activation function (Guo et al., 2017).

Specifically, the network architecture of the RNN implemented is a LSTM layer followed by the final softmax layer. Moreover, we performed a hyperparameter optimization based on the random search optimization provided by the Deep Learning for Java library. For all of the layers, the hyperparameter optimization resulted in an initial learning rate of 0.02 and the Normal distribution described in He et al. (2015) for the weight initialization. In addition, for the LSTM layer, the hyperparameter optimization resulted in a layer size of 223 and the standard sigmoid activation function as the activation function.

Since RNN is also based on supervised learning such as the TLR-FNN baseline, training and testing processes are also required. Therefore, we have performed the same encoding as the TLR-FNN baseline (See Fig. 12).

6.2.5. TLR-LtoR: Learning to Rank baseline

Taking into account this baseline, we want to determine if the better results of TLR-EltoR are due to the combination of the evolutionary algorithm and LtoR, or there is no need to combine the

Table 1

Mean Values and Standard Deviations for Precision, Recall, F-Measure, and Matthews Correlation Coefficient (MCC) for the baselines and the TLR-ELtoR approach

	Precision	Recall	F-Measure	MCC
TLR-Linguistic	37.38 \pm 16.18	48.61 \pm 19.78	40.41 \pm 16.19	0.40
TLR-IR	18.09 \pm 25.55	53.45 \pm 38.70	21.69 \pm 23.95	0.21
TLR-FNN	8.20 \pm 0.10	100 \pm 0.00	14.06 \pm 0.14	-0.84
TLR-RNN	8.37 \pm 0.09	100 \pm 0.00	14.34 \pm 0.14	-0.77
TLR-LtoR	13.01 \pm 26.08	11.85 \pm 18.24	10.27 \pm 17.57	0.07
TLR-ELtoR	59.91 \pm 33.39	78.95 \pm 15.16	62.50 \pm 27.76	0.64

two to get these results. For this purpose, this baseline is based only on LtoR, and the model fragments that are used as input for the LtoR process are generated randomly through a standard random search.

We used this algorithm as outlined in Algorithm 2 (available in Font et al., 2017). The algorithm starts with a random initial model fragment, as the best fragment. A new random model fragment is then generated and assessed using LtoR. Then, the values provided by LtoR for both fragments, the best one and the new one, are compared and the model fragment with the greatest value is selected as the best one. The search then goes back to the second step, generating and assessing a new model fragment, and this loop is repeated until a stop condition is met.

Therefore, this baseline does not take advantage of evolving model fragments to guide the exploration of the models, as our approach does thanks to the evolutionary algorithm. Since there is no a guide to explore the models, the search for the model fragment that realizes a specific requirement is performed by brute-force.

Since TLR-LtoR is also based on LtoR such as TLR-ELtoR, training and testing are also required. Therefore, the same steps that were described for TLR-ELtoR in Section 5 are applied in TLR-LtoR to encode model fragments as feature vectors, to train a classifier from the knowledge base, and to test the test cases. In addition, the technical details, such as the LtoR algorithm and the cross-validation method, are also the same ones defined in the setup of TLR-ELtoR (See Section 6.1.1). However, the stop condition is different from TLR-ELtoR in order so that the comparison between them is fair.

The stop condition in TLR-ELtoR was set up to perform 200 iterations of the evolutionary algorithm, where each iteration evaluated 120 model fragments. Therefore, for each test case, the approach evaluated a total of 24,000 model fragments. However, TLR-LtoR approach only evaluates one model fragment for each iteration, so the stop condition was set up to perform 24,000 iterations in order to evaluate the same number of model fragments.

6.3. Case study

The case study where we applied our approach was CAF, a worldwide provider of railway solutions. Their trains can be found all over the world and in different forms (regular trains, subway, light rail, monorail, etc.). A train unit is furnished with multiple pieces of equipment in its vehicles and cabins. These pieces of equipment are often designed and manufactured by different providers, and their aim is to carry out specific tasks for the train. Some examples of these devices are: the traction equipment, the compressors that feed the brakes, the pantograph that harvests power from the overhead wires, and the circuit breaker that isolates or connects the electrical circuits of the train. The control software of the train unit is in charge of making all of the equipment cooperate in order to achieve the train functionality, while guaranteeing compliance with the specific regulations of each country. The following video illustrates the CAF models: [youtube.com/watch?v=Ypci2evEQB8](https://www.youtube.com/watch?v=Ypci2evEQB8)

Our evaluation includes 20 test cases, which are composed of a requirement, a product model, a knowledge base, and an ontology. A detailed description of each of them and how they are used in our approach TLR-ELtoR is provided below:

- The **requirements** have about 25 words. In TLR-ELtoR, these requirements and the models are used to generate the testing sets.
- The **models** have about 650 elements. In TLR-ELtoR, model fragments are generated by the evolutionary algorithm through genetic operations from these models. Then, these model fragments and the requirements compose the testing set.
- The **knowledge base** includes 103 samples. Specifically, each of these samples contains a requirement, a model fragment that has about 15 elements, and an assessment. In TLR-ELtoR, the entire knowledge base composes the training set that is used to train the classifier. In the end, the number of samples in the knowledge base allows the generation of a suitable, but not perfect, classifier (Beleites et al., 2013). Fig. 14 shows the distribution of scores in the knowledge base, taking into account how many samples of the knowledge base are in each scores range.
- The **ontology** contains a total of 54 elements between concepts and relations. In TLR-ELtoR, the ontology is used to encode the knowledge base in the training process, and to encode the requirements and the models in the testing process.

For each test case, we followed the experimental setup described in Fig. 7. Each test case was run 30 times. As suggested by Arcuri and Fraser (2013), given the stochastic nature of the TLR-ELtoR approach, several repetitions are needed to obtain reliable results. Finally, the results were evaluated and compared to the oracle. The oracle contains the **approved traceability**, which consists of a set of model fragments, where each model fragment contains the model elements that are required by the requirement. In other words, the oracle contains the solutions for each test case, so the oracle had 20 model fragments, one for each test case.

6.4. Implementation details

We used the Eclipse Modeling Framework to manipulate the models and CVL to manage the model fragments. For the development of the TLR-Linguistic baseline, the Stanford POS Tagger (sta, 2017) was utilized. The LSI technique used within the TLR-IR baseline was implemented using the Efficient Java Matrix Library (EJML Abeles, 2017). The neural networks in TLR-FNN and TLR-RNN were developed and tuned by means of the Deep Learning for Java library (Team et al., 2016). The genetic operations were built upon the Watchmaker Framework for Evolutionary Computation (Dyer, 2016). Finally, RankBoost was implemented using the RankLib library (Dang, 2013).

6.5. Results

In Table 1, we outline the results, which are aggregated for each of the baselines and for our approach. Each row shows the Precision, Recall, F-measure, and MCC obtained through each technique.

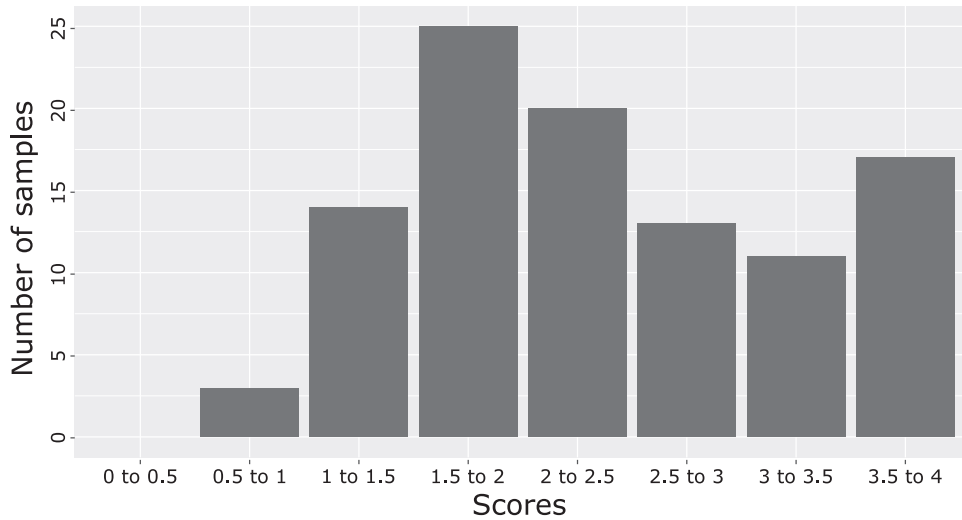


Fig. 14. Distribution of scores in the knowledge base

As the table shows, TLR-ELtoR achieves the best results for most performance indicators, providing a mean precision value of 59.91%, a recall value of 78.95%, a combined F-measure value of 62.50%, and a MCC value of 0.64. In contrast, the TLR-Linguistic baseline, the TLR-IR baseline, and the TLR-LtoR baseline present worse results in all of the measurements: the TLR-Linguistic baseline attains 37.38% precision, 48.61% recall, 40.41% F-measure, and 0.40 MCC; the TLR-IR baseline achieves 18.09% precision, 53.45% recall, 21.69% F-measure, and 0.21 MCC; and the TLR-LtoR baseline attains 13.01% precision, 11.85% recall, 10.27% F-measure, and 0.07 MCC. On the other hand, both the TLR-FNN baseline and the TLR-RNN baseline achieve the best results for recall, but they present the worst results for the rest of the indicators: the TLR-FNN attains 8.20% precision, 100% recall, 14.06% F-measure, and -0.84 MCC; and the TLR-RNN baseline achieves 8.37% precision, 100% recall, 14.34% F-measure, and -0.77 MCC.

7. Statistical analysis

To properly compare the different configurations, the data resulting from the empirical analysis was analyzed using statistical methods.

7.1. Statistical significance

A statistical test must be run to assess whether there is enough empirical evidence to claim that there is a difference between two approaches (e.g., A is better than B). To achieve this, two hypotheses are defined: the null hypothesis H_0 , and the alternative hypothesis H_1 . The null hypothesis H_0 is typically defined to state that there is no difference between the approaches, whereas the alternative hypothesis H_1 states that the configurations differ. In such a case, a statistical test aims to verify whether the null hypothesis H_0 should be rejected.

Statistical tests provide a probability value, p -Value. The p -Value obtains values between 0 and 1. The lower the p -Value of a test, the more likely that the null hypothesis is false. It is accepted by the research community that a p -Value under 0.05 is statistically significant (Arcuri and Briand, 2014), and so the hypothesis H_0 can be considered false.

The test carried out depends on the properties of the data. Since our data does not follow a normal distribution in general, our analysis required the use of nonparametric techniques. There are several tests for analyzing this kind of data; however, the

Table 2

Quade test statistic and p -Values

	Recall	Precision
p-Value	2.20×10^{-16}	1.7×10^{-10}
Statistic	35.27	14.41

Table 3

Holm's Post Hoc p -Values

	Recall	Precision
TLR-Linguistic vs TLR-IR	0.49	6.1×10^{-04}
TLR-Linguistic vs TLR-FNN	1.0×10^{-07}	3.4×10^{-07}
TLR-Linguistic vs TLR-RNN	1.0×10^{-07}	3.4×10^{-07}
TLR-Linguistic vs TLR-LtoR	1.9×10^{-05}	2.5×10^{-03}
TLR-Linguistic vs TLR-ELtoR	3.3×10^{-06}	8.4×10^{-03}
TLR-IR vs TLR-FNN	4.2×10^{-06}	0.04
TLR-IR vs TLR-RNN	4.2×10^{-06}	0.04
TLR-IR vs TLR-LtoR	1.1×10^{-03}	0.27
TLR-IR vs TLR-ELtoR	0.041	7.3×10^{-05}
TLR-FNN vs TLR-RNN	0.0	1.1×10^{-03}
TLR-FNN vs TLR-LtoR	2.9×10^{-08}	0.97
TLR-FNN vs TLR-ELtoR	8.3×10^{-06}	1.9×10^{-07}
TLR-RNN vs TLR-LtoR	2.9×10^{-08}	0.97
TLR-RNN vs TLR-ELtoR	8.3×10^{-06}	1.9×10^{-07}
TLR-LtoR vs TLR-ELtoR	2.7×10^{-07}	3.8×10^{-06}

Quade test is the most powerful one when working with real data (García et al., 2010). In addition, according to Conover (1999), the Quade test is the one that has shown the best results for a low number of approaches (no more than 4 or 5 approaches).

Table 2 shows the Quade test statistic and p -Values for recall and precision. Since the p -Values are smaller than 0.05, we rejected the null hypothesis. Consequently, we can state that there are differences among the five approaches.

Nevertheless, with the Quade test, we cannot answer the following question: *Which of the approaches gives the best performance?* In this case, the performance of each approach should be individually compared with all of the other alternatives. In order to do this, we performed an additional post hoc analysis. This kind of analysis performs a pair-wise comparison among the results of each approach, determining whether statistically significant differences exist among the results of a specific pair of approaches.

Table 3 shows the p -Values of Holm's post hoc analysis for each specific pair of approaches. Almost of all the p -Values shown in this table are smaller than 0.05, except for some cases:

Table 4
 \hat{A}_{12} statistic for each pair of approaches

	Recall	Precision
TLR-Linguistic vs TLR-IR	0.45	0.81
TLR-Linguistic vs TLR-FNN	0.0	0.93
TLR-Linguistic vs TLR-RNN	0.0	0.93
TLR-Linguistic vs TLR-LtoR	0.91	0.86
TLR-Linguistic vs TLR-ELtoR	0.17	0.30
TLR-IR vs TLR-FNN	0.13	0.55
TLR-IR vs TLR-RNN	0.13	0.54
TLR-IR vs TLR-LtoR	0.78	0.65
TLR-IR vs TLR-ELtoR	0.32	0.15
TLR-FNN vs TLR-RNN	0.5	0.48
TLR-FNN vs TLR-LtoR	1	0.67
TLR-FNN vs TLR-ELtoR	0.85	0.07
TLR-RNN vs TLR-LtoR	1	0.68
TLR-RNN vs TLR-ELtoR	0.85	0.07
TLR-LtoR vs TLR-ELtoR	0.03	0.11

the recall comparison between TLR-Linguistic and TLR-IR, the recall comparison between TLR-FNN and TLR-RNN, the precision comparison between TLR-IR and TLR-LtoR, the precision comparison between TLR-FNN and TLR-LtoR, and the precision comparison between TLR-RNN and TLR-LtoR. Therefore, significant differences for one of the performance measurements were obtained in all of the comparisons.

7.2. Effect size

Statistically significant differences can be obtained even if they are so small as to be of no practical value [Arcuri and Briand \(2014\)](#). It is then important to assess whether an approach is statistically better than another and to assess the magnitude of the improvement. *Effect size* measures are needed to analyze this.

For a non-parametric effect size measure, we used Vargha and Delaney's \hat{A}_{12} ([Vargha and Delaney, 2000](#)). \hat{A}_{12} measures the probability that running one approach yields higher values than running another approach. If the two approaches are equivalent, then \hat{A}_{12} will be 0.5.

For example, $\hat{A}_{12} = 0.7$ means that we would obtain better results in 70% of the runs with the first of the pair of approaches that have been compared, and $\hat{A}_{12} = 0.3$ means that we would obtain better results in 70% of the runs with the second of the pair of approaches that have been compared. Thus, we have an \hat{A}_{12} value for every pair of approaches.

[Table 4](#) shows the values of the effect size statistics between every pair of approaches.

TLR-Linguistic vs TLR-IR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-Linguistic will obtain better results in 58% of the cases for recall, while TLR-IR will obtain better precision values in 81% of the cases.

TLR-Linguistic vs TLR-FNN: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-FNN will obtain better results in 100% of the cases for recall, while TLR-Linguistic will obtain better precision values in 93% of the cases.

TLR-Linguistic vs TLR-RNN: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results in 100% of the cases for recall, while TLR-Linguistic will obtain better precision values in 93% of the cases.

TLR-Linguistic vs TLR-LtoR: The \hat{A}_{12} measure value indicates that TLR-Linguistic will obtain better results than TLR-LtoR in 91% of the cases for recall, and better precision values in 86% of the cases.

TLR-Linguistic vs TLR-ELtoR: The \hat{A}_{12} measure value indicates that TLR-ELtoR will obtain better results than TLR-Linguistic in 83% of the cases for recall, and better precision values in 70% of the cases.

TLR-IR vs TLR-FNN: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results in 87% of the cases for recall, while TLR-IR will obtain better precision values in 55% of the cases.

TLR-IR vs TLR-RNN: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results in 87% of the cases for recall, while TLR-IR will obtain better precision values in 54% of the cases.

TLR-IR vs TLR-LtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-IR will obtain better results than TLR-LtoR in 78% of the cases for recall, and better precision values in 65% of the cases.

TLR-IR vs TLR-ELtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-ELtoR will obtain better results than TLR-IR in 68% of the cases for recall, and better precision values in 85% of the cases.

TLR-FNN vs TLR-RNN: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results in 50% of the cases for recall, and better precision values in 52% of the cases.

TLR-FNN vs TLR-LtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-FNN will obtain better results than TLR-LtoR in 100% of the cases for recall, and better precision values in 68% of the cases.

TLR-FNN vs TLR-ELtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-FNN will obtain better results in 85% of the cases for recall, while TLR-ELtoR will obtain better precision values in 93% of the cases.

TLR-RNN vs TLR-LtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results than TLR-LtoR in 100% of the cases for recall, and better precision values in 68% of the cases.

TLR-RNN vs TLR-ELtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-RNN will obtain better results in 85% of the cases for recall, while TLR-ELtoR will obtain better precision values in 93% of the cases.

Table 5
Required artifacts for each approach

Artifacts	Approaches					
	TLR-Linguistic	TLR-IR	TLR-FNN	TLR-RNN	TLR-LtoR	TLR-ELtoR
Models	X	X	X	X	X	X
Requirements	X	X	X	X	X	X
Knowledge Base			X	X	X	X
Ontology			X	X	X	X
Rules	X					

Table 6
Artifacts whose properties have an impact on the results

Artifacts	Properties				
	Homogeneity	Completeness	Heterogeneity	Size	Volume
Models	X			X	
Requirements	X	X			
Knowledge Base		X	X	X	
Ontology		X		X	
Rules		X			X

TLR-LtoR vs TLR-ELtoR: The \hat{A}_{12} measure value indicates that, of the two approaches, TLR-ELtoR will obtain better results than TLR-LtoR in 97% of the cases for recall, and better precision values in 89% of the cases.

The obtained \hat{A}_{12} values show that TLR-ELtoR is superior to all of the baselines for precision. Moreover, TLR-ELtoR is also superior to TLR-Linguistic, TLR-IR, and TLR-LtoR on recall, meaning that TLR-ELtoR will obtain better results than these three approaches in most of the cases. Overall, these measurements confirm that, for recall and precision, TLR-ELtoR outperforms the baselines (TLR-IR and TLR-Linguistic) that obtain the best results for TLR between requirements and models. Moreover, these measurements confirm that, for precision, TLR-ELtoR outperforms the ML baselines (TLR-FNN and TLR-RNN) that have successfully been applied recently in TLR. Finally, these measurements confirm that, for recall and precision, TLR-ELtoR outperforms the baseline (TLR-LtoR) that explores the search space by means of brute-force.

8. Discussion

In this section, we discuss what prerequisites are needed by each approach, what properties affect the results and limit the approaches. We also discuss why TLR-ELtoR is less sensitive to tacit knowledge and vocabulary mismatch than the baselines. These advantages lead to the better results of TLR-ELtoR.

8.1. Prerequisites and properties

Both our approach and the approaches in the baselines need some prerequisites to be applied. If one of their prerequisites is not satisfied, the approach would not be used in that domain. Table 5 shows what artifacts are needed to apply each approach.

Table 5 shows that all of the approaches need models and requirements. Specifically, the models where requirements have to be located must conform to MOF (the OMG metalanguage for defining modeling languages) and that requirements must be provided using natural language. Moreover, all of the approaches that are based on Machine Learning (TLR-ELtoR, TLR-LtoR, TLR-FNN, and TLR-RNN) need a knowledge base to train and an ontology to encode the models and requirements. Specifically, the knowledge base must be composed of a set of feature vectors with

the format described in Joachims (1999), and the ontology must contain a set of concepts and the relations with each other (See Fig. 6). Finally, the TLR-Linguistic approach needs rules to identify relations between model elements and requirement words. The rules have to be defined following the guides and examples in Spanoudakis et al. (2004).

Therefore, even though the training in TLR-ELtoR is beneficial in avoiding to a large extent issues such as tacit knowledge and vocabulary mismatch, it is necessary to have access to a knowledge base and an ontology to perform the training. In industrial domains, especially long-living ones, where requirements and models have been stored for years, a knowledge base may be easily available. Also, thanks to the wide experience of the employees in companies of this kind, the main concepts and relations could be identified by experts in the domain. However, in other scenarios, such as when only the first product has been developed, TLR-ELtoR cannot be applied.

In addition, even though we had all the necessary artifacts to apply our approach, the results may not be as good as possible. In fact, some properties of the artifacts have an impact on the results. For example, if there is not enough information in the knowledge base, TLR-ELtoR would not train properly, so the results would be worse than expected. Table 6 shows the properties that we have identified in this work and that had an impact on the obtained results.

The following paragraphs provide more details about the properties identified in Table 6:

- **Models** may be developed by several engineers and at different times, so the terms used to describe model elements may be different (e.g., pantograph and panto are two different terms used in our models to refer to the same concept: *pantograph*). Therefore, the first model property that affects to the results is homogeneity. The second one is the size, which has an impact on the result based on the understanding of the models (Störrle, 2014).
- **Requirements** may be defined by different engineers and at different times, so the homogeneity of the requirements, like the homogeneity of the models, has an impact on the results. The results are also affected by the completeness of the requirements. Often, when requirements are written, part of the domain knowledge related to the requirements is not embodied in them because tacit knowledge about the domain is assumed to be known by all of the domain experts. Therefore, the requirements are more or less complete in accordance with how many assumptions are made by the engineers. In the end, requirements may lose part of the information that is required because of these assumptions.
- **Knowledge base** contains the information necessary to train the classifier, so this information must be enough to train the classifier. If the knowledge base only contains the information to recover the traces between one requirement and one model, the classifier may not learn how to recover the traces for other requirements or models. Therefore, including heterogeneity samples of traces in the knowledge base provides more

complete information for the training. In addition, some Machine Learning techniques require a larger knowledge base than others to provide suitable results, so the technique must be selected based on the available knowledge base.

- **Ontology** is composed of the main concepts and relations of a domain. Therefore, if a relevant concept or relation is not present in the ontology, the encoding for the fitness function will not take it into account and the training may be incomplete, leading to worse results. For this reason, the first property to keep in mind for the ontology is completeness. Moreover, if the ontology contains unnecessary concepts or relations, the number of features for the encoding would be greater and a great number of features in the training step leads to overfitting. Therefore, we must also take into account the size of the ontology.
- **Rules** are defined by humans through the manual comparison of the models and requirements. Therefore, the completeness of the rules depends on how well engineers understand the models and requirements and how complex these models and requirements are. The volume of rules also affects the results. If only one rule is defined, the approach only recovers one type of model element, so the approach may need several rules. However, a large number of rules does not guarantee the best results. Therefore, both completeness and volume must be taken into account.

8.2. Advantages of TLR-ELtoR

This section discusses why TLR-ELtoR achieves better results than the baselines regarding three aspects: tacit knowledge, vocabulary mismatch, and available documentation.

8.2.1. Tacit knowledge

Often, when requirements are written, part of the domain knowledge related to the requirements is not embodied in them. The tacit knowledge about the domain is assumed to be known by all of the domain experts, so it is never formalized in writing. This behavior has been reported in previous works (Rus and Lindvall, 2002; Stone and Sawyer, 2006). For example, given the requirement: *At all stations, the doors are automatically opened*, the engineers understand that the doors have to be opened in all of the stations, without being requested by a passenger. However, this requirement also embodies tacit knowledge that is not written but is obvious to the domain engineers: *The train has doors on both sides, but only the doors on the side of the platform will be opened, while the doors on the side of the tracks will remain closed, and all of the doors on one side will be opened, except the driver's door in the cabin*.

The tacit knowledge is not reflected in the text of the requirements. This tacit knowledge is shared among the engineers that write the requirements and the engineers that read the requirements. Therefore, both the text of the requirements and tacit knowledge are used to build the models. As a result, the model contains elements that are related to text of the requirement, but the model also contains elements that are related to the tacit knowledge. However, since part of the knowledge is not reflected in the text of the requirement, recovering the most relevant model fragment for a requirement is complex.

Both TLR-IR and TLR-Linguistic depend, to a large extent, on the text of the requirement. TLR-IR evaluates the similarity between the requirement and the model fragment according to the co-occurrences of terms between the two. TLR-Linguistic evaluates the similarity between the requirement and the model fragment according to patterns that relate the terms in the requirement with the elements in the model fragment. In both cases, the lack of terms that is caused by the tacit knowledge makes it impossible

to locate the elements from the model that are relevant to the requirement.

In contrast, TLR-ELtoR is less sensitive to tacit knowledge due to training. In the training, the requirements of the knowledge base are linked to the model fragments of the knowledge base. Even though the text of requirements is inaccurate due to tacit knowledge, the linked model fragments are complete. Consequently, the classifier is not only trained from the text of the requirements, but also from the elements of the model fragments. Therefore, the classifier learns that certain elements of models are relevant to certain requirements even though these elements are not described properly in the text of the requirements. As TLR-ELtoR depends, to a lesser extent, on the text of the requirement than TLR-IR and TLR-Linguistic, when the requirements have a lack of terms due to tacit knowledge, the results that are obtained through TLR-ELtoR are better than the results obtained through TLR-IR and TLR-Linguistic.

8.2.2. Vocabulary mismatch

Vocabulary mismatch is caused by the use of different terms to reference the same concept in the requirement and the model. In industrial environments, sometimes the engineer who is in charge of writing the requirement is not the same engineer assigned to building the model. Moreover, both the requirement and the model may be manipulated by different engineers.

Even though TLR-IR, TLR-Linguistic, and even TLR-ELtoR, may use Natural Language Processing (NLP) to homogenize the terms between requirements and models, vocabulary mismatch continues to be an issue that must be taken into account. Since the in-house terms that are used in a specific domain or company are not known synonyms, these in-house terms may not be included in NLP, causing vocabulary mismatch. For example, the terms *PLC* and *system* may be recognized as synonyms, but the terms *PLC* and *COSMOS* are definitely not known to be synonyms because *COSMOS* is an in-house term that is used exclusively by our industrial partner to refer to *PLC*.

As in the tacit knowledge issue, TLR-IR and TLR-Linguistic are seriously affected by vocabulary mismatch because both of them depend, to a large extent, on the text of the requirements. If the terms that are used in the requirements and the terms that are used in the models are not known synonyms, they cannot be related, and therefore the requirement cannot be correctly related to the elements of the model. Therefore, the lack of awareness that is caused by vocabulary mismatch makes it impossible to locate the elements from the model that are relevant to the requirement.

In contrast, TLR-ELtoR is less sensitive to vocabulary mismatch for the same reason described for the tacit knowledge issue. The evaluation of TLR-ELtoR depends on the information provided by training. If the information that is extracted through the training indicates that a term of the requirement is related to a term of an element in the model, the classifier learns that both terms are related to each other even when they are not considered synonyms. Therefore, TLR-ELtoR depends, to a lesser extent, on the synonyms than TLR-IR and TLR-Linguistic, which leads to our approach having better results than the baselines.

8.2.3. Available documentation

Since TLR-FNN and TLR-RNN are trained using the same knowledge base than TLR-ELtoR, they should also be less sensitive to tacit knowledge and vocabulary mismatch. However, our knowledge base may be unsuitable for properly training a Neural Network. For example, in Guo et al. (2017), the training set is composed of 45% of the 769,366 artifacts, so this training set contains about 423,151 feature vectors. However, our training set is composed of the encoding of the knowledge base that has 103 samples whose model fragments have around 15 elements. Therefore, since

the ending is performed at the model-element level, the training set contains about 1545 (103 x 15) feature vectors.

Some works analyze the impact of the number of samples on the performance of the neural networks. The authors in [Mather and Koch \(2011\)](#); [Piper \(1992\)](#) suggest the use of a minimum of 1030p samples for training, where p is the number of features vectors used. However, this rule is often universally enforced in remote sensing without questioning its relevance to the complexity of the specific problem [Van Niel et al. \(2005\)](#). In fact, in some domains, the best result are obtained with 2p or 4p samples for training ([Walczak and Cerpa, 1999](#); [Van Niel et al., 2005](#)). Therefore, a small knowledge base may be insufficient and a large knowledge base may introduce noise.

On the other hand, the knowledge base may also be affected by the vaporization problem ([van der Ven et al., 2006](#)). In fact, some industrial companies do not store enough information to create a knowledge base with the necessary completeness and size. However, these domains also need to recover the traceability links, and our approach can be successfully used even if the knowledge base is small, as our evaluation proves.

Since TLR-LtoR is based on LtoR as TLR-ELtoR also is, we might expect that TLR-LtoR will not have the problems described. However, TLR-LtoR obtained the worst results because the search space was too big, so the exploration of this search space randomly required many more iterations. Therefore, by evaluating the same number of model fragments using the two approaches, the TLR-ELtoR obtained the best results thanks to the combination of the LtoR, which provides a successful evaluation of the model fragments, and the evolutionary algorithm, which allows the search space to be explored in an effective way.

9. Threats to validity

In this section, we use the classification of threats to validity of [Wohlin et al. \(2012\)](#) to acknowledge the limitations of our approach.

Construct validity: This aspect of validity reflects the extent to which the operational measures that are studied represent what the researchers have in mind. To minimize this risk, our evaluation is performed using four measures: precision, recall, F-measure, and MCC. These measures are widely accepted in the software engineering research community.

Internal Validity: This aspect of validity is of concern when causal relations are examined. There is a risk that the factor being investigated may be affected by other neglected factors. RankBoost tends to overfit when the knowledge base is not large enough and there are many encoding features [Wolf and Martin \(2005\)](#). The number of samples in our knowledge base may look small; however, this threat has been reduced because our approach uses only 54 encoding features, which is a small number in machine learning applications ([Zhou and Feng, 2017](#); [Wang et al., 2014](#)).

External Validity: This aspect of validity is concerned with to what extent it is possible to generalize the findings, and to what extent the findings are of relevance for other cases. Both requirements and models are frequently leveraged to specify all kinds of different software. The requisites for applying our approach are that the set of models where the requirements must be located

conform to MOF (the OMG metalanguage for defining modeling languages), and that the requirements must be provided using natural language. Our experiment does not rely on the specific conditions of our domain. Nevertheless, the experiment and its results should be replicated in other domains before assuring their generalization.

Reliability: This aspect is concerned with to what extent the data and the analysis are dependent on the specific researchers. To reduce this threat, the knowledge base, the requirements descriptions, and the product models were provided by our industrial partner.

10. Related work

Works that are related to our research are mainly found within the knowledge area of Traceability Link Recovery. In a more general fashion, works in the knowledge area of Feature Location can be relevant for our research as well. Feature Location takes a query as input and returns its materialization as a result. Potentially, the techniques that are applied to locate features could also be applied to Traceability Link Recovery for requirement queries. In this section, we analyze some of the existing approaches in both areas and compare our work with these approaches.

10.1. Traceability link recovery

There are several TLR techniques in use that utilize requirement queries. Most of them deal with source code or focus on the usage and impact of the specific techniques in use. This section analyzes these kinds of works and differentiates our work from them.

Most of the existing works focus on Traceability Link Recovery between requirements and source code. CERBERUS ([Eaddy et al., 2008](#)) provides a hybrid technique that combines information retrieval, execution tracing, and prune dependency analysis allowing the tracing of requirements to source code. [Eaddy et al. \(2007\)](#) present a systematic methodology for identifying which code is related to which requirement, and a suite of metrics for quantifying the amount of crosscutting code. [Marcus and Maletic \(2003b\)](#) use LSI for recovering the traceability relations between source code and documentation (manuals, design documentation, requirement documents, test suites, etc.). [Antoniol et al. \(2002\)](#) propose a method based on information retrieval to recover traceability links between source code and free text documents, such as requirement specifications, design documents, manual pages, system development journals, error logs, and related maintenance reports. [Zisman et al. \(2003\)](#) automate the generation of traceability relations between textual requirement artifacts and object models using heuristic rules. These approaches recover the traceability between source code and requirements. In contrast, our work recovers the traceability between requirements and models.

In [Marcén et al. \(2017\)](#), the authors propose an evolutionary ontological encoding approach to enable Machine Learning techniques to be used to perform Software Engineering tasks in models. Their proof of concept consists of recovering traceability links between requirements and model fragments. However, in the real world whole models are available, rather than model fragments. Therefore, the application of their approach to a real-world problem may be impossible or hard. In contrast, our work recovers the traceability links between requirements and whole models thanks to the evolutionary algorithm that generates and maintains the population of model fragments from whole models.

Some other works focus on the impact and application of Linguistics to TLR problem resolution at several levels of abstraction. Works like [Sultanov and Hayes \(2010\)](#); [Sundaram et al. \(2010\)](#) or [Duan and Cleland-Huang \(2007\)](#), among many others, use Linguistic approaches to tackle specific TLR problems and tasks. In [Falessi et al. \(2013\)](#), the authors use Linguistic techniques to identify the equivalence between requirements, also defining and using a series of principles for evaluating their performance when identifying equivalent requirements. The authors of [Falessi et al. \(2013\)](#) conclude that, in their field, the performance of Linguistic techniques is determined by the properties of the given dataset over which they are performed. They measure the properties as a factor to adjust the Linguistic techniques accordingly and then apply their principles to an industrial case study. The work presented in [Arora et al. \(2015\)](#) uses Linguistic techniques to study how changes in requirements impact other requirements in the same specification. In their work, the authors analyze TLR between requirements and use Linguistic techniques to determine how changes in requirements must propagate.

Our work differs from [Ryan \(1993\)](#); [Sultanov and Hayes \(2010\)](#); [Sundaram et al. \(2010\)](#); [Duan and Cleland-Huang \(2007\)](#) since our approach is not based or focused on Linguistic techniques as a means of TLR analysis; instead, we use an evolutionary algorithm to perform TLR between requirements and models, using Linguistic techniques only as a baseline for our work. Moreover, we do not study how Linguistic techniques must be tweaked for specific problems as [Falessi et al. \(2013\)](#) does. In addition, in contrast to [Arora et al. \(2015\)](#), we do not tackle changes in requirements on TLR between requirements, but instead focus our work on TLR between requirements and a set of evolving models.

Some recent works focus on improving TLR results through Neural Networks. [Guo et al. \(2017\)](#) present a solution to improve the current automated techniques, which fail to understand the semantics of the software artifacts or to integrate domain knowledge into the tracing process. Therefore, they tend to deliver imprecise and inaccurate results. Specifically, they utilize Word Embedding and Recurrent Neural Network (RNN) models to generate trace links, which contain the requirements artifact semantics and the domain knowledge. [Zhao et al. \(2017\)](#) propose training deep neural networks for generating text-based knowledge in software repositories to improve the accuracy of TLR. The authors in [Antoniol et al. \(2017\)](#) present some challenges in traceability and some of their proposals consider solving these traceability issues through neural networks. In our work, we do not use Neural Networks to improve the results or to perform TLR; instead, we use them as a baseline for comparison with our main line of work, TLR between requirements and models. In addition, we do not address the traceability between source codes and requirements as most of these works do.

Finally, other works target the application of LSI to TLR tasks. [De Lucia et al. \(2004\)](#) present a Traceability Link Recovery method and tool based on LSI in the context of an artifact management system, which includes models. The work in [Eder et al. \(2015\)](#) takes in consideration the possible configurations of LSI when using the technique for TLR between requirements artifacts, namely requirements and test cases. In their work, the authors state that the configurations of LSI depend on the datasets used, and they expect to be able to determine automatically an appropriate configuration for LSI for any given dataset. In our work, we do not use LSI to perform TLR; instead, we use it as a baseline for comparison with our main line of work, TLR between requirements and models. In addition, we do not tackle different LSI configurations or how LSI configurations impact the results of TLR between requirements and models as [Eder et al. \(2015\)](#) does.

10.2. Feature location

There are several Feature Location techniques and approaches that are applied to locate features for requirement queries, and, as such, can potentially be used for TLR purposes. This section covers these sorts of works, and compares our work with those.

Typechef [Kästner et al. \(2011\)](#) provides an infrastructure to locate the code that is associated to a given feature by means of analyzing the `#ifdef` directives. Trace analysis ([Eisenberg and Volder, 2005](#)) is a run-time technique that is used to locate features. When the technique is executed, it produces traces indicating which parts of code have been executed. Some approaches related to Feature Location use LSI to extract the code that is associated to a feature ([Poshyvanyk et al., 2007](#); [Liu et al., 2007](#)). These techniques have generally been applied to search the code of a feature in a given individual product. In contrast, our approach searches for model fragments that implement a requirement.

Feature Location approaches in a product family, such as the one presented in [Xue et al. \(2012\)](#), center their efforts on finding the code that implements a feature among the different products by combining techniques such as FCA ([Ganter and Wille, 1997](#)) and LSI. In our approach, we are not interested in the code representation of a feature in the family but in locating the most relevant model fragments that implement a requirement. Other works such as [She et al. \(2011\)](#) focus on applying reverse engineering to the source code to obtain the variability model. In [Czarnecki and Wasowski \(2007\)](#), the authors use propositional logic, which describes the dependencies between features. In [Nadi et al. \(2014\)](#), the authors combine Typechef and propositional logic to extract conditions among a collection of features. These works explicitly engage the variability of products, but they do not indicate the most relevant model fragments for the development of requirements, as our work does.

In [Lapeña et al. \(2016\)](#), Lapeña et al. use Linguistic techniques in combination with an adapted two-step LSI to obtain rankings of methods for all of the requirements of a new product in a product family. The scope of our work is centered around finding model fragments that can be used to implement a specific requirement, while ([Lapeña et al., 2016](#)) focuses on finding relevant code for the implementation methods of all of the new requirements in a new product in a family.

In addition, even though we had all the necessary artifacts to apply our approach, the results may not be as good as possible. In fact, some properties of the artifacts have an impact on the results. For example, if there is not enough information in the knowledge base, TLR-ELtoR would not train properly, so the results would be worse than expected. [Table 6](#) shrequirements.

[Font et al. \(2016a\)](#) use a Single Objective Evolutionary Algorithm (SOEA) to locate features among a family of models in the form of a variation point. Their approach is refined in [Font et al. \(2016b\)](#), where the authors use a SOEA to find sets of suitable feature realizations. The authors first cluster model fragments based on their common attributes into feature realization candidates through Formal Concept Analysis, and then LSI ranks the candidates based on the similarity with the feature description. In contrast, our presented approach locates model fragments for requirements instead of variation points for features. In addition, the approaches by Font et al. use FCA and LSI, while our approach trains a LtoR classifier from legacy products to guide the evolutionary algorithm.

Several approaches rely on evolutionary Algorithms guided by LSI for Feature Location. In [Font et al. \(2017\)](#), Font et al. performed a comparison among five different evolutionary algorithms for feature location in models, showing that the best results were achieved by a hybrid between an evolutionary algorithm and a hill climbing technique. In [Cetina et al. \(2017\)](#), they explored a new di-

rection: taking advantage of already long-living software systems (designed with sustainability in mind) to address the challenge of feature location. Specifically, they used commonality and modifications fitness through model retrospectives in order to promote model fragments that have undergone less modification over time. In contrast, in our approach, the target of the evolutionary algorithm is TLR instead of Feature Location. Moreover, our algorithm is guided by LtoR instead of LSI, which according to our results is less sensitive to the tacit knowledge issues of requirements.

Some approaches rely on LtoR algorithms to locate features in the code [Ye et al. \(2014\)](#); [Le et al. \(2016\)](#). Tien-Duy et al. focus on LtoR using feature vectors that are based on likely invariants. Xin et al. focus on the terms that are defined in a vocabulary to build the feature vectors. In our approach, we use LtoR as an objective for the evolutionary algorithm and perform TLR between requirements and models instead of Feature Location in code.

Other works rely on ontologies to locate features in code. In [Wilson \(2010\)](#), a systematic approach is used to locate features by using ontology fragments. [Hayashi et al. \(2010\)](#) propose an ontology-based technique to locate features that are defined by natural language descriptions. [Ratiu and Deisenboeck \(2007\)](#) present a framework to recover the mappings between entities from an ontology and program elements. [Petrenko et al. \(2008\)](#) perform a study about the performance of programmers when they locate features by using ontology fragments. In contrast, our approach performs TLR between requirements and models, using an evolutionary algorithm that is guided by a LtoR classifier.

Finally, the works presented in [Zhang et al. \(2011, 2012\)](#); [Wille et al. \(2013\)](#) focus on the location of features in models using comparisons among models in a family of models. [Zhang et al. \(2011\)](#) propose a generic approach to locate the feature realizations by exploring the commonality and the variability of models through their automatic comparison. In [Zhang et al. \(2012\)](#), the approach is refined to reduce the manual effort required in the formalization of the feature realizations when new product models are included in a product line. In the approach presented in [Wille et al. \(2013\)](#), the variability between models is determined using an exchangeable metric, taking into account different attributes of the models. However, all of these approaches are based on the location of features through comparisons among the models, while our approach performs TLR between requirements and models. In addition, we do so by relying on an evolutionary algorithm guided by a LtoR classifier.

11. Replication of the results

The implementation for our approach is available at <http://bitbucket.org/svitusj/flame>. We have also made the dataset available in the same url. The dataset contains the requirements and the models that are used in our experiment as well as the knowledge base with the requirements and the model fragments that are used to train the classifier. The implementation for the five baselines is also available at the same location. Therefore, our public online repository contains the source code of our approach, the source code of the two baselines, and the dataset (requirements and models).

12. Conclusions

Both Evolutionary Algorithms and Learning to Rank algorithms have a wide range of successful applications, but current research efforts have so far neglected the application of the two on Traceability Link Recovery (TLR) between requirements and models. In this paper, we propose the TLR-ELtoR approach, which recovers traceability links between the requirements of a software system

and its models by leveraging the usage of an evolutionary algorithm (EA) that is guided by a Learning to Rank (LtoR) algorithm.

We evaluated our TLR-ELtoR approach in terms of precision, recall, the F-measure, and the Matthews Correlation Coefficient. To do this, we compared it to five baselines in an industrial domain (firmware of train PLCs with CAF). The first baseline is a Linguistic Rule-Based (TLR-Linguistic) approach that is based on Parts-of-Speech (POS) tagging and traceability rules. The second one is an Information Retrieval (TLR-IR) approach that is based on Latent Semantic Indexing (LSI) and Singular Value Decomposition (SVD). The third one is a Feedforward Neural Network (TLR-FNN) approach that is based on a traditional neural network structure. The fourth one is a Recurrent Neural Network (TLR-RNN) approach that is based on an extension of a Feedforward Neural Network with feedback connections to model the temporal characteristics of the problem being learned ([Engelbrecht, 2007](#)). The fifth one is a Learning to Rank (TLR-LtoR) approach based on ranking Machine Learning algorithms of the same name. We report our evaluation, including: experimental setup, results, statistical analysis, and threats to validity.

The results show that the application of an evolutionary algorithm guided using the LtoR algorithm by means of TLR-ELtoR pays off for TLR. The results also show that our approach can be applied in real-world environments. The statistical analysis of the results assesses the level of the improvement that our approach offers. Moreover, the discussion shows how our approach is limited by the available documentation and how our approach may be beneficial for dealing with issues such as tacit knowledge and vocabulary mismatch.

We acknowledge that we could have proposed other approaches as baselines. For example, instead of using LSI or SVD, we could have used LtoR or other machine learning technique to support TLR-IR in order to check the relevance of a model element with a requirement. The baselines were selected according to their performance, taking into account a classification of approaches for TLR in models, or their recent successful, taking into account recent approaches for TLR. However, a future work could consist of a deeper comparison of our approach with other approaches for TLR, where the contribution could be not only the comparison with other alternative baselines but also the discussion of the advantages and limitations of these alternatives for TLR in models.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This work has been developed with the financial support of the [Spanish State Research Agency](#) and the [Generalitat Valenciana](#) under the projects DataME TIN2016-80811-P, ALPS RTI2018-096411-B-I00, ACIF/2018/171 and PROMETEO/2018/176, and co-financed with ERDF.

References

- Abeles, P., 2017. Efficient Java Matrix Library. [Online; accessed 12-April-2017], <http://ejml.org/>.
- Affenzeller, M., Winkler, S.M., Wagner, S., Beham, A., 2009. Genetic Algorithms and Genetic Programming - Modern Concepts and Practical Applications. CRC Press.
- Antoniol, G., Canfora, G., Casazza, G., De Lucia, A., Merlo, E., 2002. Recovering traceability links between code and documentation. *IEEE Trans. Softw. Eng.* 28 (10), 970–983.
- Antoniol, G., Cleland-Huang, J., Hayes, J. H., Vierhauser, M., 2017. Grand Challenges of Traceability: The Next Ten Years. *arXiv:1710.03129*.

- Arcuri, A., Briand, L., 2014. A Hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering. *Softw. Test. Verif. Reliab.* 24 (3), 219–250.
- Arcuri, A., Fraser, G., 2013. Parameter tuning or default values? an empirical investigation in search-based software engineering. *Empir. Softw. Eng.* 18 (3), 594–623. doi:10.1007/s10664-013-9249-9.
- Arora, C., Sabetzadeh, M., Goknil, A., Briand, L.C., Zimmer, F., 2015. Change impact analysis for natural language requirements: an NLP approach. In: IEEE 23rd International Requirements Engineering Conference.
- Beleites, C., Neugebauer, U., Bocklitz, T., Krafft, C., Popp, J., 2013. Sample size planning for classification models. *Analytica Chimica Acta* 760, 25–33.
- Bengio, Y., Simard, P., Frasconi, P., 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Netw.* 5 (2), 157–166.
- Bianchini, M., Maggini, M., Jain, L.C., 2013. *Handbook on Neural Information Processing*. Springer.
- Canuto, S.D., Belém, F.M., Almeida, J.M., Gonçalves, M.A., 2013. A comparative study of learning-to-rank techniques for tag recommendation. *J. Inf. Data Manage.* 4 (3), 453.
- Cao, Z., Tian, Y., Le, T.-D.B., Lo, D., 2018. Rule-based specification mining leveraging learning to rank. *Autom. Softw. Eng.* 1–30.
- Cetina, C., Font, J., Arcega, L., Pérez, F., 2017. Improving feature location in long-living model-based product families designed with sustainability goals. *J. Syst. Softw.* 134, 261–278.
- Chandrashekar, G., Sahin, F., 2014. A survey on feature selection methods. *Comput. Electr. Eng.* 40 (1), 16–28.
- Chopra, C., Sinha, S., Jaroli, S., Shukla, A., Maheshwari, S., 2017. Recurrent neural networks with non-sequential data to predict hospital readmission of diabetic patients. In: Proceedings of the 2017 International Conference on Computational Biology and Bioinformatics. ACM, pp. 18–23.
- Conover, W., 1999. *Practical Nonparametric Statistics*, 3rd ed.. Wiley, New York, pp. 250–257.
- Czarnecki, K., Wasowski, A., 2007. Feature diagrams and logics: there and back again. In: Proceedings of the 11th International Software Product Lines Conference.
- Dang, V., 2013. The Lemur Project - Wiki - RankLib. [Online; accessed April-2017], <http://sourceforge.net/p/lemur/wiki/RankLib/>.
- Davis, L., 1991. *Handbook of Genetic Algorithms*.
- De Lucia, A., Fasano, F., Oliveto, R., Tortora, G., 2004. Enhancing an artefact management system with traceability recovery features. In: Proceedings of the 20th IEEE International Conference on Software Maintenance. IEEE, pp. 306–315.
- De Lucia, A., Fasano, F., Oliveto, R., Tortora, G., 2006. Can information retrieval techniques effectively support traceability link recovery? In: 14th IEEE International Conference on Program Comprehension. IEEE, pp. 307–316.
- Duan, C., Cleland-Huang, J., 2007. Clustering support for automated tracing. In: Proceedings of the 22nd IEEE/ACM International Conference on Automated Software Engineering.
- Dyer, D., 2016. The Watchmaker Framework for Evolutionary Computation (Evolutionary/Genetic Algorithms for Java). [Online; accessed 7-April-2016], <http://watchmaker.uncommons.org/>.
- Eaddy, M., Aho, A., Murphy, G.C., 2007. Identifying, assigning, and quantifying cross-cutting concerns. In: Proceedings of the First International Workshop on Assessment of Contemporary Modularization Techniques, p. 2.
- Eaddy, M., Aho, A.V., Antoniol, G., Guéhéneuc, Y.-G., 2008. Cerberus: tracing requirements to source code using information retrieval, dynamic analysis, and program analysis. In: ICPC 2008 conference. IEEE, pp. 53–62.
- Eder, S., Femmer, H., Hauptmann, B., Junker, M., 2015. Configuring latent semantic indexing for requirements tracing. In: Proceedings of the 2nd International Workshop on Requirements Engineering and Testing.
- Eisenberg, A.D., Volder, K.D., 2005. Dynamic feature traces: finding features in unfamiliar code. In: 21st IEEE International Conference on Software Maintenance.
- Engelbrecht, A.P., 2007. *Computational Intelligence: An Introduction*, 2nd ed. Wiley Publishing.
- Fallessi, D., Cantone, G., Canfora, G., 2013. Empirical principles and an industrial case study in retrieving equivalent requirements via natural language processing techniques. *Trans. Softw. Eng.* 39 (1).
- Font, J., Arcega, L., Haugen, Ø., Cetina, C., 2016. Feature location in model-based software product lines through a genetic algorithm. In: Proceedings of the 15th International Conference on Software Reuse: Bridging with Social-Awareness, pp. 39–54.
- Font, J., Arcega, L., Haugen, Ø., Cetina, C., 2016. Feature location in models through a genetic algorithm driven by information retrieval techniques. In: Proceedings of the ACM/IEEE 19th International Conference on Model Driven Engineering Languages and Systems. ACM, pp. 272–282.
- Font, J., Arcega, L., Haugen, Ø., Cetina, C., 2017. Achieving feature location in families of models through the use of search-based software engineering. *IEEE Trans. Evol. Comput.*
- Frakes, W. B., Baeza-Yates, R., 1992. *Information Retrieval: Data Structures and Algorithms*.
- Freund, Y., Iyer, R., Schapire, R.E., Singer, Y., 2003. An efficient boosting algorithm for combining preferences. *J. Mach. Learn. Res.* 4 (Nov), 933–969.
- Ganter, B., Wille, R., 1997. *Formal Concept Analysis: Mathematical Foundations*, 1st ed. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- García, S., Fernández, A., Luengo, J., Herrera, F., 2010. Advanced nonparametric tests for multiple comparisons in the design of experiments in computational intelligence and Data Mining: Experimental Analysis of Power. *Inf. Sci.* 180 (10), 2044–2064.
- Ghazarian, A., 2010. A research agenda for software reliability. IEEE Reliability Society 2009 Annual Technology Report.
- Gotel, O.C., Finkelstein, C., 1994. An analysis of the requirements traceability problem. In: Proceedings of the First International Conference on Requirements Engineering. IEEE, pp. 94–101.
- Guo, J., Cheng, J., Cleland-Huang, J., 2017. Semantically enhanced software traceability using deep learning techniques. In: Software Engineering (ICSE), 2017 IEEE/ACM 39th International Conference on. IEEE, pp. 3–14.
- Haugen, Ø., Møller-Pedersen, B., Oldevik, J., Olsen, G.K., Svendsen, A., 2008. Adding standardized variability to domain specific languages. In: Proceedings of the 12th International Software Product Lines Conference, pp. 139–148.
- Hayashi, S., Yoshikawa, T., Saeki, M., 2010. Sentence-to-Code traceability recovery with domain ontologies. In: 2010 Asia Pacific Software Engineering Conference, pp. 385–394.
- Haykin, S., 1994. *Neural Networks: a Comprehensive Foundation*. Prentice Hall PTR.
- He, K., Zhang, X., Ren, S., Sun, J., 2015. Delving deep into rectifiers: surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision, pp. 1026–1034.
- Hirzel, A.H., Le Lay, G., Helfer, V., Randin, C., Guisan, A., 2006. Evaluating the ability of habitat suitability models to predict species presences. *Ecol. Model.* 199 (2), 142–152.
- Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Comput.* 9 (8), 1735–1780.
- Hornik, K., Stinchcombe, M., White, H., 1989. Multilayer feedforward networks are universal approximators. *Neural Netw.* 2 (5), 359–366.
- Joachims, T., 1999. SvmLight: Support vector machine. SVM-Light Support Vector Machine <http://svmlight.joachims.org/>, University of Dortmund 19 (4).
- Kästner, C., Giarrusso, P.G., Rendel, T., Erdweg, S., Ostermann, K., Berger, T., 2011. Variability-aware parsing in the presence of lexical macros and conditional compilation. In: Proceedings of the 26th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications.
- Kırac, M.F., Aktemur, B., Sözer, H., 2018. VISOR: a fast image processing pipeline with scaling and translation invariance for test oracle automation of visual output systems. *J. Syst. Softw.* 136, 266–277.
- Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S., 2017. Self-normalizing neural networks. In: Advances in neural information processing systems, pp. 971–980.
- Landauer, T.K., Foltz, P.W., Laham, D., 1998. An introduction to latent semantic analysis. *Discourse Process.* 25 (2–3), 259–284.
- Lapeña, R., Ballarín, M., Cetina, C., 2016. Towards clone-and-own support: locating relevant methods in legacy products. In: Proceedings of the 20th International Conference on Software Product Lines.
- Lapeña, R., Font, J., Pastor, Ó., Cetina, C., 2017. Analyzing the impact of natural language processing over feature location in models. *ACM SIGPLAN Not.* 52 (12), 63–76.
- Le, T.-D.B., Lo, D., Le Goues, C., Grunske, L., 2016. A learning-to-rank based fault localization approach using likely invariants. In: Proceedings of the 25th International Symposium on Software Testing and Analysis. ACM, pp. 177–188.
- Leech, G., Garside, R., Bryant, M., 1994. CLAWS4: the tagging of the British National Corpus. In: Proceedings of the 15th Conference on Computational Linguistics - Volume 1. Association for Computational Linguistics, pp. 622–628.
- Liu, D., Marcus, A., Poshvanyk, D., Rajlich, V., 2007. Feature location via information retrieval based filtering of a single scenario execution trace. In: Proceedings of the 22nd IEEE/ACM International Conference on Automated Software Engineering. ACM, pp. 234–243.
- Lu, Z.-J., Xiang, Q., Wu, Y.-m., Gu, J., 2015. Application of support vector machine and genetic algorithm optimization for quality prediction within complex industrial process. In: Industrial Informatics (INDIN), 2015 IEEE 13th International Conference on. IEEE, pp. 98–103.
- Lucia, A.D., Fasano, F., Oliveto, R., Tortora, G., 2007. Recovering traceability links in software artifact management systems using information retrieval methods. *ACM Trans. Softw. Eng. Methodol.* (TOSEM) 16 (4), 13.
- Lucia, D., et al., 2000. Information retrieval models for recovering traceability links between code and documentation. In: Proceedings of the International Conference on Software Maintenance. IEEE, pp. 40–49.
- Manning, C.D., Raghavan, P., Schütze, H., et al., 2008. *Introduction to information retrieval*, 1. Cambridge University Press.
- Mao, J., Xu, W., Yang, Y., Wang, J., Huang, Z., Yuille, A., 2014. Deep captioning with multimodal recurrent neural networks (M-RNN). arXiv:1412.6632.
- Marcén, A.C., Pérez, F., Cetina, C., 2017. Ontological evolutionary encoding to bridge machine learning and conceptual models: approach and industrial evaluation. In: International Conference on Conceptual Modeling. Springer, pp. 491–505.
- Marcus, A., Maletic, J.L., 2003. Recovering documentation-to-source-code traceability links using latent semantic indexing. In: Proceedings of the 25th International Conference on Software Engineering. IEEE, pp. 125–135.
- Marcus, A., Maletic, J.L., 2003. Recovering documentation-to-source-code traceability links using latent semantic indexing. In: Proceedings of the 25th International Conference on Software Engineering. IEEE, pp. 125–135.
- Marcus, A., Sergeyev, A., Rajlich, V., Maletic, J., 2004. An information retrieval approach to concept location in source code. In: Proceedings of the 11th Working Conference on Reverse Engineering, pp. 214–223. doi:10.1109/WCRE.2004.10.
- Mather, P.M., Koch, M., 2011. *Computer Processing of Remotely-Sensed Images: an Introduction*. John Wiley & Sons.
- Meziane, F., Athanasakis, N., Ananiadou, S., 2008. Generating natural language specifications from UML class diagrams. *Requirements Eng.* 13 (1), 1–18.

- Nadi, S., Berger, T., Kästner, C., Czarnecki, K., 2014. Mining configuration constraints: static analyses and empirical results. In: 36th International Conference on Software Engineering.
- Oliveto, R., Gethers, M., Poshyvanyk, D., De Lucia, A., 2010. On the equivalence of information retrieval methods for automated traceability link recovery. In: 18th International Conference on Program Comprehension. IEEE, pp. 68–71.
- Parizi, R.M., Lee, S.P., Dabbagh, M., 2014. Achievements and challenges in state-of-the-art software traceability between test and code artifacts. IEEE Trans. Reliab. 63 (4), 913–926.
- Petrenko, M., Rajlich, V., Vanciu, R., 2008. Partial domain comprehension in software evolution and maintenance. In: 16th IEEE International Conference on Program Comprehension, pp. 13–22.
- Piper, J., 1992. Variability and bias in experimentally measured classifier error rates. Pattern Recognit. Lett. 13 (10), 685–692.
- Poshyvanyk, D., Guéhéneuc, Y., Marcus, A., Antoniol, G., Rajlich, V., 2007. Feature location using probabilistic ranking of methods based on execution scenarios and information retrieval. IEEE Trans. Softw. Eng. 33 (6).
- Ratiu, D., Deissenboeck, F., 2007. From reality to programs and (not quite) back again. In: 15th IEEE International Conference on Program Comprehension. IEEE, pp. 91–102.
- Refaeilzadeh, P., Tang, L., Liu, H., 2009. Cross-validation. In: Encyclopedia of database systems. Springer, pp. 532–538.
- Rempel, P., Mäder, P., 2017. Preventing defects: the impact of requirements traceability completeness on software quality. IEEE Trans. Softw. Eng. 43 (8), 777–797.
- Rummler, A., Grammel, B., Pohl, C., 2007. Improving traceability in model-driven development of business applications. In: ECMDA Traceability Workshop (ECMDA-TW), pp. 7–15.
- Rus, I., Lindvall, M., 2002. Knowledge management in software engineering. IEEE Softw. 19 (3), 26.
- Ryan, K., 1993. The role of natural language in requirements engineering. In: Proceedings of IEEE International Symposium on Requirements Engineering.
- Salman, H.E., Seriai, A., Dony, C., 2014. Feature location in a collection of product variants: combining information retrieval and hierarchical clustering. In: The 26th International Conference on Software Engineering and Knowledge Engineering, pp. 426–430.
- Salton, G., McGill, M. J., 1986. Introduction to modern information retrieval.
- Sayyad, A.S., Ingram, J., Menzies, T., Ammar, H., 2013. Scalable product line configuration: a straw to break the Camel's Back. In: IEEE/ACM 28th International Conference on Automated Software Engineering (ASE), pp. 465–474. doi:10.1109/ASE.2013.6693104.
- Shabtai, A., Moskovitch, R., Elovici, Y., Glezer, C., 2009. Detection of malicious code by applying machine learning classifiers on static features: a state-of-the-art survey. information security technical report 14 (1), 16–29.
- She, S., Lotufo, R., Berger, T., Wasowski, A., Czarnecki, K., 2011. Reverse engineering feature models. In: Proceedings of the 33rd International Conference on Software Engineering.
- Song, Q., Jia, Z., Shepperd, M., Ying, S., Liu, J., 2011. A general software defect-prone prediction framework. IEEE Trans. Softw. Eng. 37 (3), 356–370.
- Spanoudakis, G., Zisman, A., 2005. Software Traceability: a Roadmap. Handbook Softw. Eng. Knowl.Eng. 3, 395–428.
- Spanoudakis, G., Zisman, A., Pérez-Minana, E., Krause, P., 2004. Rule-based generation of requirements traceability relations. J. Syst. Softw. 72 (2), 105–127.
- Stone, A., Sawyer, P., 2006. Using pre-requirements tracing to investigate requirements based on tacit knowledge. In: ICSoft (1), pp. 139–144.
- Större, B., 2014. On the impact of layout quality to understanding UML diagrams: size matters. In: International Conference on Model Driven Engineering Languages and Systems. Springer, pp. 518–534.
- Sultanov, H., Hayes, J.H., 2010. Application of swarm techniques to requirements engineering: requirements tracing. In: 18th IEEE International Requirements Engineering Conference.
- Sundaram, S.K., Hayes, J.H., Dekhtyar, A., Holbrook, E.A., 2010. Assessing traceability of software engineering artifacts. Requirements Eng. 15 (3).
- Team, D., et al., 2016. Deeplearning4j: Open-source Distributed Deep Learning for the JVM. 2. Apache Software Foundation License.
- The Stanford Natural Language Processing Group (2017). <https://nlp.stanford.edu/software/tagger.shtml>. [Online; accessed 18-May-2017].
- Van Niel, T.G., McVicar, T.R., Datt, B., 2005. On the relationship between training sample size and data dimensionality: Monte Carlo analysis of broadband multi-temporal classification. Remote Sens. Environ. 98 (4), 468–480.
- Vargha, A., Delaney, H.D., 2000. A critique and improvement of the CL common language effect size statistics of McGraw and Wong. J. Educ. Behav. Stat. 25 (2), 101–132.
- van der Ven, J.S., Jansen, A.G., Nijhuis, J.A., Bosch, J., 2006. Design decisions: the bridge between rationale and architecture. In: Rationale management in software engineering. Springer, pp. 329–348.
- Walczak, S., Cerpa, N., 1999. Heuristic principles for the design of artificial neural networks. Inf. Softw. Technol. 41 (2), 107–117.
- Wang, J., Zhao, P., Hoi, S.C., Jin, R., 2014. Online feature selection and its applications. IEEE Trans. Knowl. Data Eng. 26 (3), 698–710.
- Watkins, R., Neal, M., 1994. Why and how of requirements tracing. IEEE Softw. 11 (4), 104–106.
- Wille, D., Holthusen, S., Schulze, S., Schaefer, I., 2013. Interface variability in family model mining. In: 17th International Software Product Line Conference.
- Wilson, L.A., 2010. Using ontology fragments in concept Location. In: IEEE International Conference on Software Maintenance, pp. 1–2.
- Winkler, S., Pilgrim, J., 2010. A survey of traceability in requirements engineering and model-driven development. Softw. Syst. Model. (SoSyM) 9 (4), 529–565.
- Wohlin, C., Runeson, P., Höst, M., Ohlsson, M.C., Regnell, B., Wesslén, A., 2012. Experimentation in Software Engineering. Springer Science & Business Media.
- Wolf, L., Martin, I., 2005. Robust boosting for learning from few examples. In: Computer Vision and Pattern Recognition, 1. IEEE, pp. 359–364.
- Xue, Y., Xing, Z., Jarzabek, S., 2012. Feature location in a collection of product variants. In: 19th Working Conference on Reverse Engineering.
- Ye, X., Bunesco, R., Liu, C., 2014. Learning to rank relevant files for bug reports using domain knowledge. In: Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, pp. 689–699.
- Zhang, X., Haugen, Ø., Møller-Pedersen, B., 2011. Model comparison to synthesize a model-driven software product line. In: 15th International Software Product Line Conference (SPLC). IEEE, pp. 90–99.
- Zhang, X., Haugen, Ø., Møller-Pedersen, B., 2012. Augmenting product lines. In: 19th Asia-Pacific Software Engineering Conference.
- Zhang, Z., Chen, L., Tian, P., Su, J., 2017. Source localization in an ocean waveguide using supervised machine learning. Computing 11, 5.
- Zhao, Y., Zaman, T.S., Yu, T., Hayes, J.H., 2017. Using Deep learning to improve the accuracy of requirements to code traceability. In: Grand Challenges of Traceability: The Next Ten Years, p. 22.
- Zhou, Z.-H., Feng, J., 2017. Deep Forest: Towards an Alternative to Deep Neural Networks. arXiv:1702.08835.
- Zisman, A., Spanoudakis, G., Pérez-Miñana, E., Krause, P., 2003. Tracing software requirements artifacts. In: Software Engineering Research and Practice, pp. 448–455.

Ana Cristina Marcén is a PhD student in Computer Science at UPV. Her research lines include model-driven development, machine learning, variability modeling, and feature location. Marcén holds a Master's Degree in Advanced Software Technologies for Mobile Devices from Universidad San Jorge (Spain), and she has received a predoctoral grant from Generalitat Valenciana ACIF/2018/171. Contact her at acmarcen@usj.es.

Raúl Lapeña is a PhD student in Computer Science at UPV. His research lines include model-driven development, latent semantic indexing, variability modeling, and feature location. Lapeña holds a Master's Degree in Advanced Software Technologies for Mobile Devices from Universidad San Jorge (Spain). Contact him at rlapena@usj.es.

Oscar Pastor is Full Professor and Director of the Research Center on “Métodos de Producción de Software (PROS)” at the Universidad Politécnica de Valencia, Spain. He received his Ph.D. in 1992. He was a researcher at HP Labs, Bristol, UK. He has published more than three hundred research papers in conference proceedings, journals and books, received numerous research grants from public institutions and private industry, and been keynote speaker at several conferences and workshops. Member of the SC of conferences as ER, CAiSE, ESEM, ICWE, CIBSE or RCIS, his research activities focus on conceptual modelling, web engineering, requirements engineering, information systems, and model-based software production. He created the object-oriented, formal specification language OASIS and the corresponding software production method OO-METHOD. He led the research and development underlying CARE Technologies that was formed in 1996. CARE Technologies has created an advanced MDA-based Conceptual Model Compiler called IntegraNova, a tool that produces a final software product starting from a conceptual schema that represents system requirements. He is currently leading a multidisciplinary project linking Information Systems and Bioinformatics notions, oriented to designing and implementing tools for Conceptual Modelling-based interpretation of the Human Genome information.

Carlos Cetina is an associate professor with San Jorge University and the Head of the SVIT Research Group. He received a PhD in computer science from the Technical University of Valencia. His research focuses on software product lines and model-driven development. His research results have reshaped software development in world-leader industries from heterogeneous domains ranging from induction hob firmware to train control and management systems. More information about his background can be found at his website: <http://carloscetina.com>.