

Evaluating few shot and Contrastive learning Methods for Code Clone Detection*

Mohamad Khajezade
University of British Columbia
Kelowna BC, Canada
khajezad@mail.ubc.ca

Fatemeh H. Fard
University of British Columbia
Kelowna BC, Canada
fatemeh.fard@ubc.ca

Mohamed S. Shehata
University of British Columbia
Kelowna BC, Canada
mohamed.sami.shehata@ubc.ca

ABSTRACT

Context: Code Clone Detection (CCD) is a software engineering task that is used for plagiarism detection, code search, and code comprehension. Recently, deep learning-based models have achieved an F1 score (a metric used to assess classifiers) of ~95% on the CodeXGLUE benchmark. These models require many training data, mainly fine-tuned on Java or C++ datasets. However, no previous study evaluates the generalizability of these models where a limited amount of annotated data is available.

Objective: The main objective of this research is to assess the ability of the CCD models as well as few shot learning algorithms for unseen programming problems and new languages (i.e., the model is not trained on these problems/languages).

Method: We assess the generalizability of the state of the art models for CCD in few shot settings (i.e., only a few samples are available for fine-tuning) by setting three scenarios: i) unseen problems, ii) unseen languages, iii) combination of new languages and new problems. We choose three datasets of BigCloneBench, POJ-104, and CodeNet and Java, C++, and Ruby languages. Then, we employ Model Agnostic Meta-learning (MAML), where the model learns a meta-learner capable of extracting transferable knowledge from the train set; so that the model can be fine-tuned using a few samples. Finally, we combine contrastive learning with MAML to further study whether it can improve the results of MAML.

KEYWORDS

Contrastive Learning, few shot learning, Code Clone Detection

ACM Reference Format:

Mohamad Khajezade, Fatemeh H. Fard, and Mohamed S. Shehata. 2022. Evaluating few shot and Contrastive learning Methods for Code Clone Detection. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Code clones are defined as identical code fragments in different software systems [5]. Software developers tend to copy and paste existing functions and libraries instead of developing every part of

their system from scratch [30]. Accordingly, code clones emerge automatically in most software development projects [29], and can cause buggy code to propagate through the whole project. Code Clone Detection (CCD) is an essential approach in developing and maintaining code [16]. Detecting code clones is useful for different applications including detecting library candidates, code comprehension, finding harmful software, and plagiarism detection [1]. There are a lot of studies that develop models to detect code clones [38–40]. Some works only detect code clones in the same language [20]. A more recent stream of studies develops models for cross-language code clone detection (CLCCD) [3, 23], where the model intends to detect code snippets with the same functionality in different programming languages.

Two gaps exist in the CCD studies. First, the generalizability of the models to unseen scenarios is unknown. Although code clone detection methods achieve an F1 score (a score used for classification task) of approximately 96 percent on the BigCloneBench data [39] and MAP@R score (a metric used for information retrieval-based CCD task) of 88 percent on POJ-014 [39], which are widely used benchmarks for code clone detection, the ability of the current state of the art models for CCD to be generalized to unseen *programming languages* and unseen *programming problems* is rarely studied (we refer to ‘programming problems’ as ‘problems’ for simplicity). For example, can a model that is trained to detect code clones in Java, achieve the same performance for detecting code clones in Python? Here, unseen programming problems refer to different classes of problems. For instance, a model trained to detect clones for a sorting problem might fail to perform well in a domain that includes string matching.

It is worth mentioning that although the CLCCD studies detect clones in different languages, the models are trained and tested on a combined set of programming languages [3, 23], meaning that if we intend to use the same model to detect the clones in a different programming language, other than the ones the model was trained on, the process of building the dataset for training and testing the model should be done from scratch, thus, reducing the ability of the model to be generalized to *unseen* programming languages.

The generalizability of the CCD models are rarely studied. If a model exists to detect clones in one language and detecting clones in a new language is required, the model should be retrained again using a *large* number of training data in the new language, but preparing such data is expensive and time consuming. Data science, where the developers are looking for the solution to the same problem in different languages [21, 23], or cross platform software development, where the companies make the same application for different platforms are among the applications of a generalizable model [23].

*This study was accepted at the MSR 2022 Registered Reports Track

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Additionally, if the code fragments are from a different problem sets, the ability of the models to detect them is shown to be downgraded significantly [46]. For example, Yu et al. [46] present TBCCD for code clone detection, which achieves an F1 score of 99 percent on a binary classification version of POJ-104 data, but they show that their model is susceptible in the face of unseen problems. The performance of their model drops to 49 percent on average when trained and tested on disjoint sets of problems (i.e. tested on unseen problems), but the authors do not provide any solutions.

Second, training CCD models requires a lot of training data, which is costly to build. Lei et al. report that 84 percent of CCD studies mainly use BigCloneBench which has over six million code pairs [38] and POJ-104 which has 52,000 code samples [17, 22]. However, many other programming languages have restricted labeled data, and generating annotated CCD data for these languages requires extensive time and effort [48]. Moreover, while there exists many submissions for some programming problems, which can be used to learn their semantics for CCD, there are problems for which only a few solutions exist, making the CCD data for training the model in that domain hard.

Given the restricted labeled data for other programming languages and the costs of curating new datasets and training the models on new datasets, there is a need for models that perform well) where there is little training data available, and 2) are generalizable to new problems sets and new programming languages. In this study, we intend to evaluate the ability of the models in three scenarios: i) the model is applied on new problem sets, ii) the model has not seen the programming language during training, or iii) a combination of both. Note that the second scenario is closely related to CLCCD, where the model is only trained on one language and is tested on another language, but different in the sense that the current CLCCD models are trained on a combination of the languages that require clone detection¹. We conduct our study in a more rigorous way and evaluate the models in cases where **only a few examples are available for training the models**, i.e., only 5, 10, or 15 samples instead of thousands or millions of samples, addressing the costs associated with curating high quality datasets and retraining the models. Additionally, to address the need for models that are able to perform well given a few labeled data, we investigate two techniques for the detection of code clone to assess if the performance of the models can be increased when only a few training records are available: few shot learning and contrastive learning integrated with the few shot algorithm.

Few shot learning is a specific approach developed for transferring the knowledge that a model has learned by training on a high resource domain/task (i.e. there exists a lot of labeled data) to a low resource domain/task (i.e. when limited labeled data is available) [48]. In few shot learning algorithms, a meta learner is trained using a target task, for which there is plenty of labeled data. This meta-learner aims to learn how to fine-tune the model using a *handful* of examples, without overfitting to those few observations. These few examples, which form a *support set*, belong to the downstream domain, for which there exists limited annotated data. In other words,

¹Thus, we do not focus on CLCCD and our focus is only of CCD in this work. It is worth mentioning that applying the same techniques can be interesting to assess the ability of the CLCCD models, but this is out of the scope of our work, due to the high number of required experiments and limited computational resources.

the only labeled data in the downstream task are the examples in the support set. Then, the fine-tuned model is used to predict the label of another set in the downstream task called the query set [42]. In this regard, note that the label set in training is disjoint with labeled samples (called 'support set') and query examples (known as 'query set') in the downstream task, for which the number of labeled data is scarce. Few shot learning models are widely studied in computer vision [31, 33, 36] and there are few works conducted on using few shot learning for natural language processing [4, 45]. However, **no previous study has been conducted on applying few shot learning algorithms for code clone detection.**

Recently, some works demonstrated the success of contrastive learning in Natural Language Processing (NLP) [37] or multimodal models like Learning transferable visual models from natural language supervision [27], which aims to gain supervision from natural language to classify images. Inspired by the achievement of contrastive learning models in NLP, some works have studied the contrastive learning objective to represent source code and for code clone detection [12, 24, 39]. In this respect, ContraCode [12] uses a self-supervised method based on contrastive learning to learn the functionality of code fragments. They report that ContraCode can improve the results of RoBERTa [19] for code clone detection by 2% to 46 % according to different experimental settings on JavaScript code. Similarly, SyncoBERT [39] is another contrastive learning approach that outperforms state of the art models for CCD. In this model, multimodal training is presented to capture the structural properties of the code snippet as well as its semantic meaning. A contrastive learning approach is used to learn the semantic similarity between modalities (i.e. code tokens, Abstract Syntax Tree graph, and identifiers) by computing the mutual information between them. Although contrastive learning has shown to benefit code clone detection, **there is no study on how these models perform for unseen programming languages and unseen problems, especially in the case of being fine-tuned only using a few examples of the language or problem of interest.**

In this study, we intend to investigate to what extent the current CCD models perform when trained on a few samples of a language or problem set that they are not trained on. We will then study a few shot learning technique for code clone detection, and finally we assess whether the of combination of few shot learning with contrastive learning can boost the performance of these models. The main novelty of this work is on evaluating code clone detection models in the context of few shot learning, for both unseen problems and unseen programming languages. The other contribution is the integration of contrastive learning with a few shot algorithm and assessing whether it can improve the performance of the CCD models. The results can shed lights on the generalizability of the clone detection models for languages in which the number of available labeled data is limited.

The rest of this report is organized as follows. In Section 2 we explain the necessary concepts of few shot and contrastive learning. Then, we mention the research questions and the two tasks of CCD in Sections 3 and 4, followed by explanation of the datasets used and our execution plan in Sections 5 and 6, respectively. Section 7 discusses the threats to validity and section ?? concludes the paper.

2 BACKGROUND

2.1 Few Shot Learning

Few shot learning is an area of machine learning where the model is trained using a restricted number of labeled data, sometimes as much as one sample in a downstream task. In this regard, few shot learning algorithms and models learn a meta learner so that the meta learner can compare the labeled samples with query examples, for which the model should predict the related class [41]. There are three different datasets in a few shot learning problem: train set, support set, and test set (also known as query set). While the support set and test set share the same label space, the label space of the train set is disjoint with both the support set and test set. In this regard, if the support set consists of K samples per each unique class and there exist C classes in the support set, then this problem is called a k -shot C -way [18]. It is possible to learn a classifier using only labeled samples in the support set. However, as the number of labeled data is scarce, the model would be overfitted to a few examples in the support set. Thus, few shot learning algorithms learn a meta-learner capable of extracting transferable knowledge from the train set so that the model can be successfully fine-tuned using labeled data in the support set and performs more accurately on the test set [18]. Accordingly, to better train the meta learner on the train set, an episode approach similar to [36] is utilized in most studies to mimic the few shot learning setting for meta-learning. In this sense, an episode is created in each training iteration by sampling C classes in the train set. Then, K -samples are chosen randomly for each class to serve as the support set. Moreover, a query set is formed using the remaining samples. The model is then trained using the support set in each episode, and the losses are computed using the query set to perform learning. We follow this episode-based approach for code clone detection.

2.2 Contrastive Learning

In contrastive learning, instead of training a model to predict a ground truth label, the model aims to produce a representation, which is similar for positive samples and is different for negative pairs [15]. Contrastive learning can be used for both supervised and unsupervised learning tasks. Contrastive learning is a powerful self-supervised approach for unsupervised learning tasks. This method enables models to learn the representations without needing annotated data. Accordingly, a strategy to generate positive and negative samples to learn the contrastive learning objective should be employed. For instance, in ContraCode, a compiler to compiler translator generates different versions of the same code as positive samples [12]. Thus, different versions of different codes are considered negative examples in this method. In a contrastive supervised learning approach, on the other hand, the objective is defined such that the samples that belong to the same class have similar embeddings. Accordingly, positive samples belong to the same label space, and negative samples have different classes.

3 RESEARCH QUESTIONS

RQ1: *What is the performance of the current state of the art models for code clone detection when they are fine-tuned on a few examples of the down stream task?*

In this research question, we investigate the performance of the current models for code clone detection for unseen classes of problems and unseen languages. Following [46], this research defines unseen classes as unseen programming problems. For instance, if the model is trained on sorting algorithms, then solutions for string matching problems are considered an unseen class. For the unseen languages, we consider C++, Java, and Ruby, as C++ and Java are widely studied [10, 39, 46] and Ruby is a low resource language (i.e. a limited number of annotated data exists for Ruby) and it has not been studied previously for CCD.²

Moreover, Ruby is different from C++ and Java in terms of syntactical structure, and the models cannot benefit from the similarity of the given input languages, thus their ability to detect code clones for languages that are structurally different can be assessed. Additionally, Ruby has less number of submissions in CodeNet (see Section 5) compared to the dominated languages, thus can be considered as low resource language. Other programming languages are not considered in our study due to the high number of experiments and limitations in computational resources.

RQ2: *What is the performance of a few shot learning method for code clone detection?* After investigating the performance of the current state-of-the-art models for code clone detection in the above mentioned scenarios, we investigate the applicability of few-shot learning models in improving the performance of the models for code clone detection. In this research question we study the performance of a popular few shot learning algorithm, Model Agnostic Meta-learning [9] for code clone detection.

RQ3: *Can we improve the performance of the few shot learning models by using a contrastive learning model as the baseline for training?*

Contrastive learning has recently been developed to represent code and has been shown to have the state of the art performance for code-related tasks, including code clone detection [12, 39], as for example, in the ContraCode study, the model learns the functionality of the code fragments, which is important in detecting code clones. So, here, we investigate whether using a contrastive learning code representation approach can increase the performance of the models in the few shot setting when integrated in the few shot learning algorithm.

4 CODE CLONE DETECTION TASKS

In this research, we study two tasks of code clone detection:

Binary Classification Task: In this task, the label of code fragments C_i and C_j is 1 if they are regarded as code clones and 0 otherwise. Accordingly, n code samples in the training set can be denoted by $T = \{(C_i, C_j, y_{i,j}) | i, j \in n, i \neq j\}$. A code clone detection algorithm based on deep learning aims to find a function Φ to compare the representation of two code fragments and assign a label to each pair based on their similarity. In most deep learning models, the cosine similarity $s_{i,j}$ is used to compare the embeddings of code samples. The equation to calculate this metric for two code fragments is as follows:

²Lei et al. report that 84 percent of CCD studies focus on Java and C and the other 16 percent mainly work on Python, C#, C++, and Go [17].

Table 1: Statistics of BigCloneBench from [34]

	Number of Examples
Train	901,028
Dev	415,416
Test	415,416

Table 2: Statistics of POJ-104 from [22]

	Number of Problems	Number of Examples
Train	64	32,000
Dev	16	8,000
Test	24	12,000

$$s_{i,j} = \frac{\Phi(C_i) \cdot \Phi(C_j)}{\|\Phi(C_i)\| \|\Phi(C_j)\|} \quad (1)$$

where $s_{i,j} \in [-1, 1]$, $\phi(c_i)$ is the embedding of C_i , and $\phi(C_j)$ is the embedding of C_j . The algorithm compares this similarity with a threshold δ to map this metric into a label between 0 and 1.

Retrieval-based Code Clone Detection: We follow the definition provided in previous works [20, 39] for this task. This task compares a query code snippet with k code samples and returns top R code fragments with the most semantic similarity. Thus, the training set is represented as follow:

$$(C_i, \{C_j, \dots, C_m\}, \{C_l, \dots, C_p\} | j - m = k \text{ and } l - p = R) \quad (2)$$

where C_i is the query code, C_j, \dots, C_m are candidates compared to the query code, C_l, \dots, C_p are detected clones, k is the number of candidates, and R is the number of top most similar code fragments to C_i in terms of semantic similarity. The cosine similarity is the metric used to compare code fragments' representation in this task. As a result, after computing the cosine similarity between query code with each candidate, the algorithm outputs the top R code fragments that are more similar to the query example.

It is worth noting that retrieval-based code clone detection is very similar to code clone search. However, code clone search has different objectives and requirements. Code clone search is a branch of code clone detection that focuses on developing search engines that aim to look for clones of a code fragment (query code) in a large corpus of code snippets [13]. Studies on code clone search focus on scalability, response time, and ranking the result set [13]. However, retrieval-based code clone detection intends to compare the query code with a set of code candidates pairwise and return a fixed number of most similar matches as code clones. Accordingly, scalability, response time, and ranking of the chosen code clones are not the focus of this research. In our study, we do not investigate the code clone search and only focus on the code clone detection, either as binary classification or retrieval-based CCD.

5 DATASETS

5.1 Datasets

In this study three different datasets would be investigated.

BigCloneBench: BigCloneBench [35] is one of the main benchmarks used for code clone detection. The original BigCloneBench Consists of 6,000,000 positive samples and 260,000 negative pairs as introduced in [38]. However, in this research, we follow the version that is provided in CodeXGlue repository³ for code clone detection. This version of BigCloneBench consists of almost 1.7 million samples. Table 1 indicates the statistics related to this dataset. The only programming language that is covered in this dataset is Java. Each sample includes the source code of two functions implemented in Java. Based on whether these code samples belong to the same clone, the label of an instance is true or false. This dataset is available online in JSON format. We choose this dataset as it is widely used as a benchmark in numerous studies [8, 38, 43, 46].

POJ-104: This dataset includes 52,000 C++ samples. Table 2 shows the statistics of this dataset. The task defined on this dataset is a Retrieval-based code clone detection. The pure dataset contains C++ source codes, aiming to solve 104 different problems. Each of these problems is a directory labeled with an index to show the problem. The source codes within each directory are considered to have the same semantic similarity. This dataset has 32,000, 8,000, and 12,000, records for train, dev, and test sets, respectively. The main reason of choosing POJ-104 data in our work is that this dataset is used in previous studies for code clone detection [8, 38, 43, 46].

CodeNet: This project composes a dataset of 14 million source code examples, each aiming to solve one of the 4,000 problems in CodeNet [26]. These problems are gathered from different online resources like AIZU Online Judge and Atcoder websites. This dataset consists of 50 programming languages. However, it is dominated by Java, Python, C, and C++. CodeNet includes useful tools to filter data based on the use cases and perform transformations like converting codes to tokens, Abstract Syntax Tree (AST), and data and control graphs. The format of the raw data is similar to the one for POJ-104 and the implementations for each problem are gathered in a specific directory. The implementation that exists in the same directory can be interpreted as the codes from the same clone. Moreover, each problem is indexed as well as each implementation so that they can be tracked down in the metadata.

We developed scripts to convert this data into both binary code clone detection and Retrieval-based code clone detection. This dataset is prepared and is ready to be used. The process of preparing this dataset is explained in the next section. The main reason to use this dataset is that it includes submissions to the same problems in multiple programming languages. This makes the dataset suitable to study the performance of the models for new languages and new problems separately, each time keeping the other variable fixed. In other words, for the same set of problems, we can study the performance of models for unseen languages; and for different languages we can study the performance of the models for new problems. This dataset is a high quality data collected by IBM.

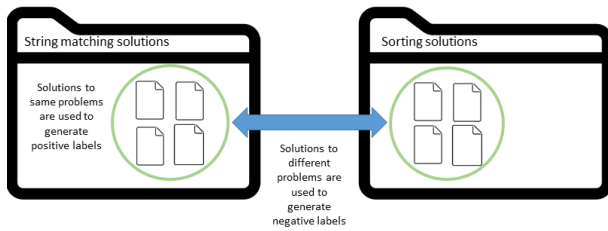


Figure 1: The procedure to extract data from CodeNet

5.2 Data Extraction

As described in the dataset subsection, the raw data for CodeNet consists of directories and files. Each directory collects all the implementations for a specific problem. An index indicates this problem. Thus, the name of the folder is the same as the index of the problem. The source codes in each directory are incorrect or correct solutions. Thus, the metadata is used to exclude the wrong samples from the final dataset. After preparing the raw data based on the task, which would be either a binary code clone detection or source code retrieval, some JSON files are generated for train, dev, and test sets.

Binary Classification Data: For generating a JSON file useful for binary clone detection, we follow the proposed schema for BigCloneBench [35]. Accordingly, the JSON file includes the following fields:

- **func1** shows the source code for the first function.
- **func2** depicts the source code for the second function.
- **index** is used to show the index of the problem.
- **id1** indicates the index of the first function.
- **id2** depicts the index of the second function.
- **label** is used to show the label of the sample, which would be either true or false.

The procedure to extract binary classification data from the CodeNet is depicted in Figure 1. Accordingly, each code sample in each directory is paired with another code snippet from the same folder for generating true labels. The idea behind creating true labels is that all the codes for the same problem share the same semantic. As a result, they belong to the same clone, which generates type 4 clones, where the code fragments are semantically the same [1]. For creating the negative labels, each source code is paired with other code snippets in other folders. In other words, implementations that belong to different problems form the negative samples. As the number of positive samples will be very large (500!), extracting this amount of data requires a huge amount of computational resources and is unnecessary. Thus, we limit the number of positive and negative samples such that the number of generated instances is the same as BigCloneBench (the version from [34]). Moreover, the procedure forces the number of true examples to be the same as negative labels, following [34].

As the few shot algorithm used in this research follows an episode-based approach to train meta-learner, having sequential labeling in the dataset overfit the meta-learner. To avoid this issue, we shuffle the dataset randomly.

³<https://github.com/microsoft/CodeXGLUE/tree/main/Code-Code/Clone-detection-BigCloneBench>

Code Retrieval Data: For extracting this dataset from CodeNet, the same raw data used for binary clone detection is considered. The JSON files generated for train, dev, and test follow the schema introduced for POJ-104 in the hugging face⁴ repository. Each JSON file includes the following fields:

- **Code** encompasses the source code for the current sample.
- **id** indicates the index of the current instance.
- **label** indicates the index of the problem that the current sample aims to solve. As mentioned, each problem has an index, and here this index is used as the label of the dataset

Extracting the dataset for code retrieval is straightforward. The algorithm traverse all directories for each problem. First, the name of the folder is assigned to the label for the current instance. Then, the algorithm processes all the submissions for the current directory and extracts their source code to be set as the code field of the data. Additionally, we will follow considerations of Ragkhitwetsagul et al. for changing the binary classification CCD data to a code retrieval data [28].

6 EXECUTION PLAN

6.1 Evaluation Metrics

Given that this research considers two different code clone detection tasks, the following metrics are used to evaluate and compare the models.

F1 Score This metric is used to evaluate classification models and is defined as follows:

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (3)$$

Precision is defined as:

$$precision = \frac{TruePositive(TP)}{TruePositives(TP) + FalsePositives(FP)} \quad (4)$$

True positives are considered as code clones (i.e. true labels) that are predicted as code clones by a code clone detection model. False positive is the number of code samples predicted as code clones, but their actual label is false.

Recall is the model's accuracy in predicting true positives and is calculated as follows:

$$Recall = \frac{TruePositives(TP)}{TruePositives(TP) + FalseNegatives(FN)} \quad (5)$$

Here, false negatives are the number of samples that are code clones but wrongly predicted with false labels.

We will report the percentage of F1 score. When comparing models, the model with higher value of F1 score is considered to perform better.

MAP@R: The mean average precision (MAP) is used for the information retrieval tasks. This metric is used for evaluating code clone detection when the model returns R similar clones between K candidates for each query code fragment. Accordingly, the precision is averaged over each class or query example. Then, the mean of average precisions for all query code fragments is computed as MAP@R. Equation 6.1 denotes the calculation for this metric:

⁴<https://huggingface.co/>

$$MAP@R = \frac{1}{R} \sum_{i=1}^R P(i) \quad (6)$$

where $P(i)$ is the precision calculated for query example i .

6.2 Baselines

BERT: BERT is a pre-trained language model that learns the contextual embedding of words from unlabeled data [7]. Accordingly, the model learns the representation of words based on their usage in different contexts. BERT trains a bidirectional text encoder by optimizing it on the Masked Language Model (MLM) and Next Sentence Prediction (NSP) tasks. For the MLM, the model first replaces 15% of the words in each sequence with the ([MASK]) token and then trains the model to predict this token based on the word's context. For the NSP task, the model takes as input two sentences and trains to predict if they are correct matches. Accordingly, 50% of sentences are true matches, and 50% of the input sentences are randomly paired for training the BERT model.

BERT is used in previous software engineering works to study software entity recognition, requirement engineering, and application review [6, 14, 32]. It is also based on the Transformer architecture [44] which is the main architecture of the pre-trained models in NLP and software engineering. Therefore, we use BERT as one of our baseline models.

RoBERTa: This model is an extension of the BERT mode. RoBERTa [19] changes the pre-training steps of BERT and achieves a better performance in various classification tasks. Accordingly, RoBERTa indicates that eliminating the NSP task from pre-training increases the performance. Moreover, in RoBERTa, longer sequences are used to optimize the MLM objective in pre-training the model. RoBERTa has been used in many state of the art approaches for code classification and representation tasks, including CodeBERT [8]. So, we use this model as a baseline in our work. For the experiments in this study, a 12 layers version of this model will be used.

CodeBERT: This model is an extension of the BERT model, which is used for understanding code-related tasks [8]. There are two versions of CodeBERT available. CodeBERT_{MLM}, uses MLM as the pre-training objective. The other version, CodeBERT_{MLM+RTD} pre-trains on MLM and Replace Token Detection (RTD) task. CodeBERT is trained on CodeSearchNet dataset [11]. This dataset consists of 6 programming languages. We use the base version of CodeBERT in our experiments. As CodeBERT is used in multiple studies in software engineering [25, 47, 49] and achieves an F1 score of 94 percent on code clone detection, we choose it as a baseline in our study.

TBCCD: TBCCD [46] is introduced as one of the state of the art models for code clone detection in the CodeXGLUE GitHub repository⁵. In this method, a siamese network is used to learn the representation of code fragments. The code fragments are first converted into AST tokens. Then, a Transformer is used to convert these tokens into an embedding. A tree-based convolutional approach processes these embeddings. The output of the convolutional layer is passed to a max pooling and a fully connected layer. Finally, the cosine similarity between code representations of this

siamese network is learned such that for non-clone pairs, the cosine similarity would be close to -1 and 1 for similar pairs.

TBCCD is chosen as one of the baselines as it is considered as the state of the art model on CodeXGLUE benchmark. It also has evaluated its model on unseen problems previously, and its setting is considered in our study.

CDLH: This model is another deep learning method [43] introduced by CodeXGLUE as a baseline for CCD. CDLH consists of two parts. Code fragments are first converted into AST tokens. Then, a tree-based LSTM model is used to create the representation of code fragments, which are passed to the second part of the model. The second part of CDLH aims to convert code representations into a binary hash code such that the hamming distance between code clones is close to each other and hamming distance for unrelated codes is far away. As this is one of the baselines for CCD as mentioned on CodeXGLUE, we choose it in our study.

ContraCode: ContraCode learns the functionality of source codes instead of capturing the syntax of the input code [12]. ContraCode employs a contrastive objective to pull the representation of positive samples together and push the negative ones apart. As the objective function, this model uses the InfoNCE, which converts the representation learning into a binary classification. A Transformer encoder is used in all experiments of ContraCode. The model is pre-trained on JavaScript source codes.

ContraCode is a contrastive learning code representation approach which has been shown to have the state of the arts results for different tasks in software engineering including code clone detection [12, 39]. Therefore, it is chosen as a baseline in our work as well as using it as the base learner in RQ3.

6.3 Experimental Setup

In this research, we aim to investigate the three RQs for the two tasks of code clone detection (i.e. binary classification and code retrieval) for three scenarios: i) unseen problems, ii) new languages, and iii) combination of unseen languages and unseen problems. Figure 2 depicts the methodology of this research. First, we evaluate the baselines models in a few shot setting as we will explain in this section (RQ1). Then, apply a widely used few shot learning approach to evaluate the models in the same settings (RQ2). Finally, we use contrastive learning in the few shot approach to evaluate the performance of the model in the few shot setting (RQ3). In the following, we first explain the three scenarios, then, we discuss the details related to each RQ. We explain the three scenarios for binary classification task first. The same settings are applied for the Retrieval-based code clone detection, but with different datasets.

Scenario I: Unseen Problems: The left part of Figure 3 illustrates the training process for this scenario. For this scenario, we choose the CodeNet data. There are over 1,492 problems for which the solutions exist in C++, Java, and Ruby. These languages are selected as C++ and Java are widely studied [17], and Ruby has not been investigated before. Moreover, Ruby has less number of submissions in CodeNet compared to the other three dominated languages.

We follow the work of Yu et al. [46] for the settings of this experiment. We only apply the experiments on 105 problems randomly selected from these 1,492 problems. The problems will be divided

⁵<https://github.com/microsoft/CodeXGLUE/tree/main/Code-Code/Clone-detection-BigCloneBench>

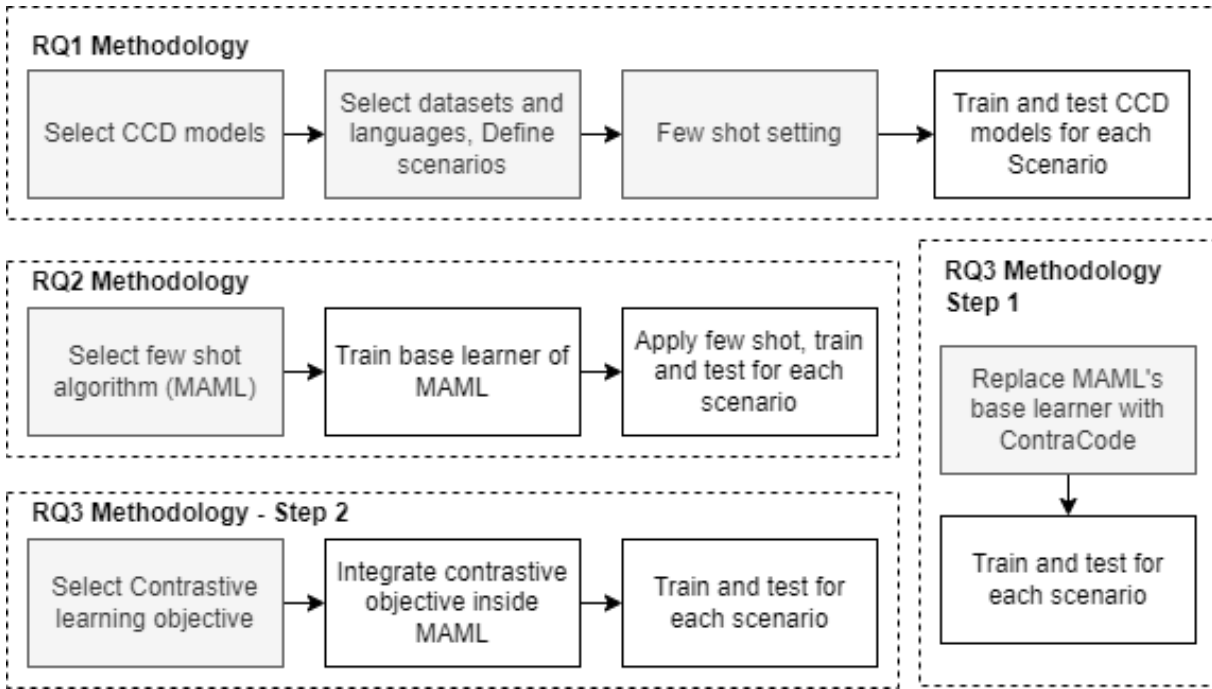


Figure 2: Big picture methodology for each research question.

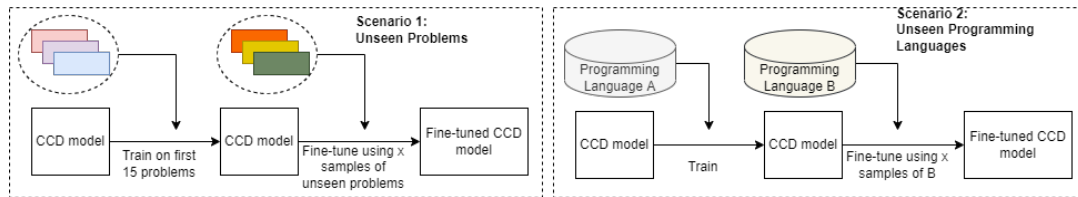


Figure 3: Left: Training process for unseen problems. The dataset is split in groups of 15 problem sets. The first 15 problems are used for training the model, and then the model is fine-tuned using few (5,10, or 15) samples of unseen problems and tested. Right: Training process for unseen programming languages. First, the model is trained on programming language A and then a few (5, 10, 15) records of a new programming language (B) is shown to the model, and the model is tested on the new programming language.

in groups of 15 classes randomly (7 groups in total). Then, we train/fine-tune each of the models separately on the first 15 problems (first group with IDs 1-15). For the few shot setting, we then fine-tune this trained model using 5 samples from the second group of problems (IDs 16-30), and then test on the rest of the problems in this group. This process is repeated for each of the problem groups with IDs 31-45, 46-60, 61-75, 76-90, and 91-105 separately, where the model is trained on the first 15 problems, fine-tuned using 5 samples of each of these groups and then tested on that group. For the few shot setting, we follow the work of Bansal et al. [4], and so repeat the same experiment once when 10 samples from each group are used for fine-tuning the models and another time with fine-tuning with 15 samples from each group. We conduct the experiments once for Java, and evaluate the best performing models based on the number of samples used for training. We refer to this number as X which will be used for the rest of experiments with

other languages, other scenarios, and Retrieval-based code clone detection in all RQs. For example, if the models perform the best using 15 samples, we will use 15 for the few shot setting and fine-tune the models with 15 samples only. We conduct the experiments once for Java, once for C++, and once for Ruby.

Note that there will be 21 results obtained for Java as we conduct these experiments separately, 7 belonging to fine-tuning with 5 samples, 7 experiments are related to fine-tuning with 10 samples, and 7 related to 15 samples. The other languages have 7 results each, for each model. In this scenario, the language is kept the same and the problems vary. We will apply the settings for sets of problems in which there are solutions in all of these four languages. The reason is to evaluate the effect of changing the problems and be able to compare the results of the languages together.

There are 499 code fragments for each problem in [46]. Here, we will also include the same number of code fragments per problem in each language.

Scenario II: Unseen Programming Languages: The right part of Figure 3 shows the training process for this scenario. For evaluating the CCD models for different programming languages, we will employ the same dataset of scenario I from CodeNet, where there are 105 problems and their solutions are available in different languages. The reason of this choice is that the problems sets for will be the same to evaluate unseen languages, which eliminates the effect of unseen problems. Here, we train the models using the Java submissions of the dataset, as it is a widely studied language for CCD. For unseen languages, C++ and Ruby are considered to evaluate the performance of the models. Then, X samples from the submission in each of the other programming languages are used to fine-tune the model. Finally, the model will be evaluated using the submissions of the selected unseen language. The fine-tuning and testing on each of the languages is done separately.

Scenario III: Unseen Problems and Languages: To investigate the mixture of unseen problems and unseen languages, BigCloneBench and the dataset of CodeNet used in scenarios I and II will be used. The baselines will be trained using the BigCloneBench dataset, and then X samples of Java submission in CodeNet will be considered to fine-tune the models. Finally, the model will be evaluated on the Java submission of CodeNet. Similarly, we will fine-tune and test the models on C++, Java, and Ruby datasets.

As there is no code retrieval version of BigCloneBench, models are trained on the code retrieval version of the POJ-104 dataset for evaluating a mixture of unseen problems and unseen programming languages. The code retrieval version of the dataset of CodeNet used in scenarios I and II will be used for fine-tuning with X samples and testing to evaluate models. All the conducted experiments are done separately for each language.

RQ1: The above mentioned scenarios and the settings are applied for RQ1. For RQ1, we evaluate the performance of all the baselines according to the few shot setting explained in each scenario and report the results.

RQ2: In RQ2, we evaluate the ability of few shot learning algorithms for code clone detection. The Model Agnostic Meta-learning Algorithm (MAML) [9] is considered to be used for this purpose. MAML is chosen as this algorithm is one of the state of the art models for few shot learning. Moreover, MAML has been applied successfully to natural language processing tasks in previous studies [4]. In MAML, the training dataset is divided into different tasks. Each task is an episode, which is formed following the approach introduced in section 2.1. Then, the best parameters for each distinct task is calculated using the following equation:

$$\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta}) \quad (7)$$

The loss of each separated task is computed. For all experiments in this paper, the cross-entropy loss is employed as the loss function of MAML following the approach introduced in [4]. Then, a query set is formed based on the episodic approach explained in section 2.1 section and is used to compute the losses for each task and are added together to optimize the meta objective. The meta optimization is calculated using the following equation:

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) \quad (8)$$

In the test time, the model trained using MAML is first fine-tuned on the downstream task, which is code clone detection in our study. Finally, the fine-tuned model predicts the label for each example in the test set.

For RQ2, BERT [7], CodeBERT [8], and RoBERTa [19] are used as the base learners of MAML separately⁶. For all the experiments using MAML, the learning rate for the outer loop is set to 5e-5 following [4]. MAML can learn the inner learning rate as one of its optimization parameters. According to [4], the internal learning rate is equal to the outer procedure in this research. Moreover, the outer epochs are set to 10, and 10 update steps are applied for each episode in training. In each episode, the support set is created by randomly choosing X samples following [31]. The remaining instances are used to form the query set to calculate the loss in each training step for each episode.

Using MAML and the base learners as mentioned, the three scenarios will be studied for this RQ, as explained previously.

RQ3: In this research question, we first use a contrastive representation of code fragments as the base learner on MAML. Second, we replace the MAML's objective with a contrastive learning function. We evaluate both settings in the three scenarios as explained previously. For the contrastive representation, the same setting as RQ2 is used to train the parameters of MAML for this research question. The based learner used in this experiment is ContraCode. ContraCode [12] is a contrastive learning model that is proposed for code representation learning. After replicating the results reported in the original paper, we conduct experiments to evaluate the performance of this contrastive learning approach for the three scenarios of unseen problems, unseen languages, and combination of new problems and languages. The same setups as explained in scenarios I-III will be applied here.

For the contrastive objective, we will replace the meta-objective function of MAML with InfoNCE [2], which is a contrastive learning objective function. For extracting positive and negative samples for training meta-objectives, we will use the submissions in the CodeNet dataset. InfoNCE converts the representation learning into a binary classification so the model predicts the positive pairs between a batch of negative samples. The following equation shows the loss function of InfoNCE:

$$\mathcal{L}_{q,k^+,k^-} = -\log \frac{\exp(q \cdot k^+ / t)}{\exp(q \cdot k^+ / t) + \sum_{k^-} \exp(q \cdot k^- / t)} \quad (9)$$

In this equation, the loss function value is low when the query q is similar to the positive key and t is the temperature hyperparameter.

7 THREATS TO VALIDITY

7.1 Internal Validity

Internal validity is associated with having unwanted results. One threat can be related to the results that will be produced in this research. To mitigate this threat, we use publicly available datasets

⁶In an initial assessment of the feasibility of using the models for this study, we tried to use TBCCD as the base learner of MAML as well. However, there is an error while training TBCCD, and therefore, we omit this model for RQ2.

for CCD and also follow a similar approach to prepare the other CCD dataset that is required for our work. The datasets are high quality datasets that are published and studied in various works. We also use common evaluation metrics employed by previous studies to report the results. We follow previous works to choose hyper-parameters of few-shot learning algorithms, and we will employ the standard procedures to train and evaluate methods including pre-trained models. All the experiments that will be executed in this research would use the same machine with the same configuration to make it easy to reproduce the results that this study reports. So we anticipate a low internal threat.

7.2 External Validity

In this study, we only consider three programming languages. So although the approach and the methodology does not differ for other programming languages, the results might not be generalizable to other programming languages. // We also only study the CCD in our work. So though the CLCCD and code clone search are closely related fields, generalizing the results to these two fields or other applications requires separate research.

7.3 Construction Validity

The metrics we choose for the evaluation of the models are the widely used and acceptable metrics for classification tasks. Both F1-score and MAP@R are previously used and reported in other CCD works. So, we do not anticipate a threat related to choosing the evaluation metrics. // Another threat could be related to choosing an appropriate few shot learning algorithm. We conducted a literature review on few shot learning algorithms, and chose the algorithm that is widely used and referenced for different applications, including natural language processing. Although few shot learning technique has not been used in software engineering for CCD, MAML is used for natural language processing previously, which is the closest application to code (compared to computer vision). Therefore, we do not anticipate a threat related to the choice of algorithm, not to mention that this is an empirical study to evaluate the application of few shot for CCD.

8 CONCLUSION

Code clone detection is the task of matching code fragments that share the same syntax or semantics and is important for software maintenance. While the recent approaches including pre-trained models, achieve the state of the art results for CCD, they are trained using large labeled datasets, and most works are only conducted on a few popular programming languages. Yet, there are programming languages that are strictly limited in labeled data, and there are programming problems, for which limited number of solutions exist. In this research, we first investigate the ability of the current CCD models when they are fine-tuned by only a few examples of the new language or problems. Then, we propose to use the Model-agnostic Meta-learning algorithm to increase the performance of these models. We employ current pre-trained models for CCD as the base learner of the MAML. In the next experiment, we replace the base learner of the MAML with ContraCode, which is a contrastive learning model for code representation learning. Finally, we propose a new contrastive objective to serve as the meta objective of MAML

for CCD, to assess whether the performance of the models can be improved when only a few training samples are available.

REFERENCES

- [1] Qurat Ul Ain, Wasi Haider Butt, Muhammad Waseem Anwar, Farooque Azam, and Bilal Maqbool. 2019. A systematic review on code clone detection. *IEEE access* 7 (2019), 86121–86144.
- [2] Laurence Aitchison. 2021. InfoNCE is a variational autoencoder. *arXiv preprint arXiv:2107.02495* (2021).
- [3] Sanjay B Ankali and Latha Parthiban. 2021. Detection and Classification of Cross-language Code Clone Types by Filtering the Nodes of ANTLR-generated Parse Tree. *International Journal of Intelligent Systems and Applications* 13, 3 (2021), 43–65.
- [4] Trapit Bansal, Rishikesh Jha, and Andrew McCallum. 2019. Learning to few-shot learn across diverse natural language classification tasks. *arXiv preprint arXiv:1911.03863* (2019).
- [5] Stefan Bellon, Rainer Koschke, Giulio Antoniol, Jens Krinke, and Ettore Merlo. 2007. Comparison and evaluation of clone detection tools. *IEEE Transactions on software engineering* 33, 9 (2007), 577–591.
- [6] Adailton Ferreira de Araújo and Ricardo Marcondes Marcacini. 2021. RE-BERT: automatic extraction of software requirements from app reviews using BERT language model. In *Proceedings of the 36th Annual ACM Symposium on Applied Computing*. 1321–1327.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [8] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155* (2020).
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*. PMLR, 1126–1135.
- [10] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366* (2020).
- [11] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436* (2019).
- [12] Paras Jain, Ajay Jain, Tianjun Zhang, Pieter Abbeel, Joseph E Gonzalez, and Ion Stoica. 2020. Contrastive code representation learning. *arXiv preprint arXiv:2007.04973* (2020).
- [13] Iman Keivanloo and Juergen Rilling. 2021. *Source Code Clone Search*. Springer Singapore, Singapore, 121–134. https://doi.org/10.1007/978-981-16-1927-4_9
- [14] MV Koroteev. 2021. BERT: a review of applications in natural language processing and understanding. *arXiv preprint arXiv:2103.11943* (2021).
- [15] Phuc H Le-Khac, Graham Healy, and Alan F Smeaton. 2020. Contrastive representation learning: A framework and review. *IEEE Access* (2020).
- [16] Maggie Lei, Hao Li, Ji Li, Namrata Aundhkar, and Dae-Kyoo Kim. 2022. Deep learning application on code clone detection: A review of current knowledge. *Journal of Systems and Software* 184 (2022), 111141.
- [17] Maggie Lei, Hao Li, Ji Li, Namrata Aundhkar, and Dae-Kyoo Kim. 2022. Deep learning application on code clone detection: A review of current knowledge. *Journal of Systems and Software* 184 (2022), 111141. <https://doi.org/10.1016/j.jss.2021.111141>
- [18] Xiaoxu Li, Zhuo Sun, Jing-Hao Xue, and Zhanyu Ma. 2021. A concise review of recent few-shot meta-learning methods. *Neurocomputing* 456 (2021), 463–468.
- [19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [20] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664* (2021).
- [21] George Varghese Mathew et al. 2021. Cross language Code Similarity and Applications in Clone Detection and Code Search. (2021).
- [22] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional neural networks over tree structures for programming language processing. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- [23] Kawser Wazed Nafi, Tonny Shekha Kar, Banani Roy, Chanchal K Roy, and Kevin A Schneider. 2019. Clcdsa: cross language code clone detection using syntactical features and api documentation. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1026–1037.
- [24] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).

- [25] Cong Pan, Minyan Lu, and Biao Xu. 2021. An empirical study on software defect prediction using codebert model. *Applied Sciences* 11, 11 (2021), 4793.
- [26] Ruchir Puri, David S Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, et al. 2021. Project CodeNet: A Large-Scale AI for Code Dataset for Learning a Diversity of Coding Tasks. *arXiv preprint arXiv:2105.12655* (2021).
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020* (2021).
- [28] Chaiyong Ragkhitwetsagul and Jens Krinke. 2019. Siamese: scalable and incremental code clone search via multiple code representations. *Empirical Software Engineering* 24, 4 (2019), 2236–2284.
- [29] Chanchal K Roy and James R Cordy. 2018. Benchmarks for software clone detection: A ten-year retrospective. In *2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 26–37.
- [30] G Shobha, Ajay Rana, Vineet Kansal, and Sarvesh Tanwar. 2021. Code Clone Detection—A Systematic Review. *Emerging Technologies in Data Mining and Information Security* (2021), 645–655.
- [31] Jake Snell, Kevin Swersky, and Richard S Zemel. 2017. Prototypical networks for few-shot learning. *arXiv preprint arXiv:1703.05175* (2017).
- [32] Chao Sun, Mingjing Tang, Li Liang, and Wei Zou. 2020. Software entity recognition method based on bert embedding. In *International Conference on Machine Learning for Cyber Security*. Springer, 33–47.
- [33] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. 2018. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1199–1208.
- [34] Jeffrey Svajlenko, Judith F Islam, Iman Keivanloo, Chanchal K Roy, and Mohamad Mamun Mia. 2014. Towards a big data curated benchmark of inter-project code clones. In *2014 IEEE International Conference on Software Maintenance and Evolution*. IEEE, 476–480.
- [35] Jeffrey Svajlenko and Chanchal K Roy. 2021. Bigclonebench. In *Code Clone Analysis*. Springer, 93–105.
- [36] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. *Advances in neural information processing systems* 29 (2016), 3630–3638.
- [37] Dong Wang, Ning Ding, Piji Li, and Hai-Tao Zheng. 2021. Cline: Contrastive learning with semantic negative examples for natural language understanding. *arXiv preprint arXiv:2107.00440* (2021).
- [38] Wenhan Wang, Ge Li, Bo Ma, Xin Xia, and Zhi Jin. 2020. Detecting code clones with graph neural network and flow-augmented abstract syntax tree. In *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 261–271.
- [39] Xin Wang, Yasheng Wang, Fei Mi, Pingyi Zhou, Yao Wan, Xiao Liu, Li Li, Hao Wu, Jin Liu, and Xin Jiang. 2021. SynCoBERT: Syntax-Guided Multi-Modal Contrastive Pre-Training for Code Representation. *arXiv preprint arXiv:2108.04556* (2021).
- [40] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859* (2021).
- [41] Yaqing Wang and Quanming Yao. 2019. Few-shot learning: A survey. (2019).
- [42] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. 2020. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)* 53, 3 (2020), 1–34.
- [43] Huihui Wei and Ming Li. 2017. Supervised Deep Features for Software Functional Clone Detection by Exploiting Lexical and Syntactical Information in Source Code. In *IJCAI*. 3034–3040.
- [44] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771* (2019).
- [45] Wenpeng Yin, Nazneen Fatema Rajani, Dragomir Radev, Richard Socher, and Caiming Xiong. 2020. Universal natural language processing with limited annotations: Try few-shot textual entailment as a start. *arXiv preprint arXiv:2010.02584* (2020).
- [46] Hao Yu, Wing Lam, Long Chen, Ge Li, Tao Xie, and Qianxiang Wang. 2019. Neural detection of semantic code clones via tree-based convolution. In *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*. IEEE, 70–80.
- [47] Xue Yuan, Guanjun Lin, Yonghang Tai, and Jun Zhang. 2022. Deep Neural Embedding for Software Vulnerability Discovery: Comparison and Optimization. *Security and Communication Networks* 2022 (2022).
- [48] Ruiheng Zhang, Shuo Yang, Qi Zhang, Lixin Xu, Yang He, and Fan Zhang. 2022. Graph-based few-shot learning with transformed feature propagation and optimal class allocation. *Neurocomputing* 470 (2022), 247–256.
- [49] Xin Zhou, DongGyun Han, and David Lo. 2021. Assessing Generalizability of CodeBERT. In *2021 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 425–436.