

A Replication Study on Predicting Metamorphic Relations at Unit Testing Level

Alejandra Duque-Torres[†], Dietmar Pfahl[†], Rudolf Ramler[‡], and Claus Klammer[‡]

[†]*Institute of Computer Science, University of Tartu, Tartu, Estonia*

E-mail: {duquet, dietmar.pfahl}@ut.ee

[‡]*Software Competence Center Hagenberg (SCCH) GmbH, Hagenberg, Austria*

E-mail: {rudolf.ramler, claus.klammer}@scch.at

Abstract—Metamorphic Testing (MT) addresses the test oracle problem by examining the relations between inputs and outputs of test executions. Such relations are known as Metamorphic Relations (MRs). In current practice, identifying and selecting suitable MRs is usually a challenging manual task, requiring a thorough grasp of the SUT and its application domain. Thus, Kanewala et al. proposed the Predicting Metamorphic Relations (PMR) approach to automatically suggest MRs from a list of six pre-defined MRs for testing newly developed methods. PMR is based on a classification model trained on features extracted from the control-flow graph (CFG) of 100 Java methods. In our replication study, we explore the generalizability of PMR. First, since not all details necessary for a replication are provided, we rebuild the entire preprocessing and training pipeline and repeat the original study in a close replication to verify the reported results and establish the basis for further experiments. Second, we perform a conceptual replication to explore the reusability of the PMR model trained on CFGs from Java methods in the first step for functionally identical methods implemented in Python and C++. Finally, we retrain the model on the CFGs from the Python and C++ methods to investigate the dependence on programming language and implementation details. We were able to successfully replicate the original study achieving comparable results for the Java methods set. However, the prediction performance of the Java-based classifiers significantly decreases when applied to functionally equivalent Python and C++ methods despite using only CFG features to abstract from language details. Since the performance improved again when the classifiers were retrained on the CFGs of the methods written in Python and C++, we conclude that the PMR approach can be generalized, but only when classifiers are developed starting from code artefacts in the used programming language.

Index Terms—Software testing, metamorphic testing, metamorphic relations, prediction modelling, replication study

I. INTRODUCTION

Metamorphic Testing (MT) is a software testing approach proposed by Chen et al. [1] to alleviate the test oracle problem. A test oracle is a mechanism for detecting whether or not the outputs of a program are correct [2], [3]. The oracle problem arises when the SUT lacks an oracle or when developing one to verify computed outputs is practically impossible [3]. MT differs from traditional testing approaches in that it examines the relations between input-output pairs of consecutive SUT executions rather than the outputs of individual SUT executions [1]. The relations between SUT inputs and outputs are known as *Metamorphic Relations* (MRs). MRs define how the outputs should vary in response to a certain change in the

input [4], [5]. In this way, testers may test the SUT indirectly by looking at whether the inputs and outputs satisfy the MRs. If an MR is violated for certain test cases, then there is a high probability that there is a fault in the SUT [4]. The most challenging task facing MT is determining suitable MRs for a particular SUT. In current practice, MRs are detected manually and require an in-depth understanding of the SUT and problem domain. As a result, the identification and selection of high-quality MRs are recognised as a big challenge.

Recently, an approach supporting unit testing at the method level was proposed by Kanewala et al., published in STVR [6] and at ISSRE [7], to “predict whether a certain method exhibits a particular MR or not”. The idea behind their Predicting Metamorphic Relations (PMR) approach is to build a model that predicts whether a method in a newly developed SUT can be tested using a specific MR. The PMR approach is based on a pre-defined set of six MRs and a classifier trained on the control-flow graph (CFG) extracted from a pool of sample methods. While the original study by Kanewala et al. showed encouraging results, it was performed on a dataset covering methods implemented in one programming language (Java) only. In order to see whether and how the PMR approach could be transferred to other programming languages, given that using CFGs supports abstracting from language and implementation details, we decided to perform further research on this.

Replication studies are required to validate experimental results achieved in prior research. They are an essential part of empirical software engineering as they prove that the observations obtained can hold (or not) under different situations. There are several forms of replication [8], [9]. An *exact* replication aims to replicate the experiments as closely as possible to the initial procedures. This demonstrates that uncontrolled random variables did not drive the initial results. In a *conceptual* replication, one or more dimensions can be modified to see how well the results hold up. It is important to note that *exact* and *conceptual* replications go by the names of *repetition* and *reproduction*, respectively, in the literature [9].

In this paper, we present a conceptual replication of the study of Kanewala et al. [6]. We follow the guidelines suggested by Carver [10] and the ACM guidelines on reproducibility (different team, different experimental setup) [11]: “The measurement can be obtained with stated precision by

a different team, a different measuring system, in a different location on multiple trials. For computational experiments, this means that an independent group can obtain the same or similar result using artefacts, which they develop completely independently”. We decided to replicate the original study in three steps. In each step, we use the same set of pre-defined MRs used in the original study.

In the first step, we rebuilt the entire pipeline for extracting CFG information from source code to training and testing classification models in order to repeat the process described in the original study by Kanewala *et al.* The authors of the original study provided the extracted CFG information for replicating their work but not the Java source code of the analyzed methods, which we retrieved from the corresponding project repositories on GitHub. We also had to develop the classification models that we then used for predicting MRs, because also they were not shared by the authors of the original paper. We conducted the first step to check whether we can re-create the classifiers with as good quality as in the original study and to prepare for the next steps by building our own classifiers based on features extracted from the CFGs derived from the source code of methods.

The next two steps of our study aim at exploring the transferability and generalizability of the PMR method. In the second step, we checked whether classifiers generated from Java code perform equally well when applied to Python and C++ code. Therefore, we created two datasets, one comprising 100 methods in Python and other one comprising 100 methods in C++. Both sets of methods implemented the exact same functionality as the 100 Java methods used in the first step. We did this to guarantee that the MRs taken from the original study would apply in the same way to the Python and C++ methods as they did in the original study using Java. To check the generalizability of the PMR method to other programming languages we then also developed individual classifiers for each programming language (Python, C++) starting out from code in the same programming language to which the classifier will be applied during evaluation.

The results of our study indicate that PMR classifiers generated based on artefacts implemented in one programming language do not perform well when applied to artefacts implemented in a different programming language, even though the classifiers are based on CFGs to abstract from programming language and implementation details, the implemented functionality is exactly identical, and the set of MRs remains unchanged. On the other hand, it seems to be possible to generalize the PMR method in the sense that it can be applied with good performance on code implemented in a different programming language as long as the PMR classifiers are redeveloped based on code implemented in the same programming language.

II. RELATED WORK

MT has been demonstrated to be an effective technique for testing in a variety of application domains, *e.g.*, autonomous driving [12], [13], cloud and networking systems [14], [15],

bioinformatic software [16], [17], scientific software [18], [19]. However, the efficacy of MT heavily relies on the specific MRs employed. Some research has been done on how to choose “good” MRs. Chen *et al.* [20] examined several MRs discovered for shortest path and critical path programs, attempting to determine MRs that are useful. Liu *et al.* [21] introduced the Composition of MRs (CMRs) technique for constructing new MRs by mixing multiple existing ones. Zhang *et al.* [22] proposed a method in which an algorithm searches for MRs expressed as linear or quadratic equations. Chen *et al.* [23] developed METRIC, a specification-based technique and related tool for identifying MRs based on the category-choice framework. They also expanded METRIC into METRIC+ by integrating the information acquired from the output domain [24]. To our knowledge, Kanewala *et al.* [6], [7], were the first to show that, in previously unseen methods, MRs can be predicted using ML techniques. They used features obtained from CFGs and a set of predefined MR to train prediction models.

III. PREDICTING METAMORPHIC RELATIONS

This study is a replication of the approach proposed by Kanewala *et al.* [6] for predicting suitable MRs for the purpose of unit testing. In this section, we first present the PMR procedure proposed by Kanewala *et al.* (Section III-A). Then, we present a detailed description of the set of pre-defined MRs used in the original and our study (Section III-B) as well as the labelled dataset used in the original study (Section III-C). Finally, we summarise the evaluation results reported in the original study (Section III-D).

A. PMR Procedure

The goal of the PMR approach is to build a model that predicts whether a method in a newly developed SUT can be automatically tested by exploiting one or more MRs contained in a pre-defined set of MRs. Figure 1 shows the PMR procedure. The PMR procedure consists of three phases. *Phase I* is responsible for creating a graph description representation derived from a method’s control-flow graph (CFG). The output of this phase is a *DOT* file. *Phase II* is in charge of feature extraction from the method’s *DOT* file. Also, each method is labelled with elements from the set of pre-defined MRs. Thus, the output of this phase is a labelled dataset. *Phase III* is in charge of training and evaluating the binary classification models that predict whether a specific MR is applicable to the unit testing of a specific method. Below we describe each phase in detail.

1) Phase I – Graph description representation: This Phase starts with the creation of the graph representation from the method’s source code. Then a labelled CFG is created by annotating each node in the CFG. The process can be split into the following two steps.

Step 1.1 – CFG generation: This step is responsible for creating the CFG from the method’s source code. To get the CFG, Kanewala *et al.* use *Soot* [25]. *Soot* generates CFG representations in *Simple* format, a typed 3-address intermediate

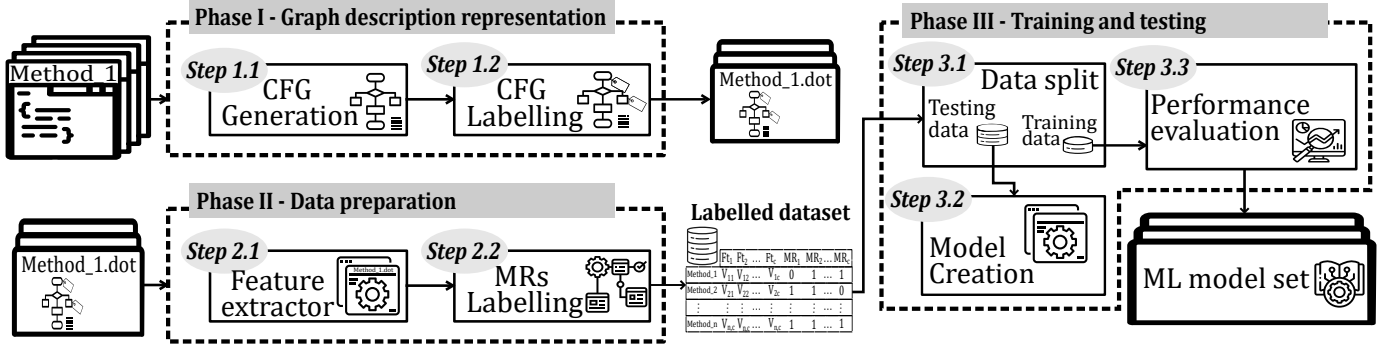


Fig. 1: PMR procedure

representation, where each CFG node represents an atomic operation [25]. The left-hand side of Figure 2 shows the CFG representation of the Algorithm 1 using the *soot* framework. The numbering of the nodes has been done manually, *i.e.*, not with the framework, and serves here only to facilitate a better understanding of the next phase and its steps.

Algorithm 1 Average of an integer array

```

1: function AVG(int input[ ])
2:   double sum = 0;
3:   double average = 0;
4:   for (int i = 0; input.length; i++) do
5:     sum += input[i];
6:   average = sum/input.length;
7:   return average

```

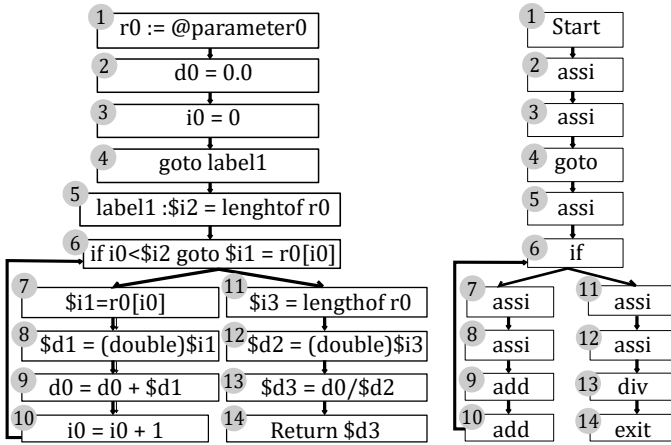


Fig. 2: On the left, CFG representation of Algorithm 1 using the *soot* framework; on the right, its CFG with annotations

Step 1.2 – CFG labelling: In this step a simplified version of the CFG is created by replacing the specific, code-related information of each node in the CFG by a more general annotation describing the specific operations and conditional jumps in the original code. The right-hand side of Figure 2 shows the CFG representation with the node annotations of the Algorithm 1. Table I shows examples of annotations that

are assigned depending on the node operation. The annotations follow the graph description language, *i.e.*, the DOT format.

2) **Phase II – Data preparation:** This phase is in charge of extracting a set of features from the annotated CFGs, *i.e.*, from the Phase I output. Also, to each method’s annotated CFG zero to six pre-defined MRs are assigned, depending on their suitability. Like Phase I, Phase II consists of two steps.

Step 2.1 – Feature extraction: Kanewala *et al.* propose two approaches for extracting features from CFG representations, features based on nodes and paths, and features based on graph similarity measures. In the former, the node features (hereafter simply NF) follows the form $NO_n - d_{in} - d_{out}$, where NO_n stands for Node Operation of node n , and d_{in} and d_{out} stand for *in-degree* and *out-degree*, respectively. The number of a specific NF type, *i.e.*, $NO_n - d_{in} - d_{out}$, is tallied and used as the NF value. As an example of NF, let us consider the annotated CFG of Algorithm 1. As the right-hand side of Figure 2 shows, the annotated CFG of Algorithm 1 has fourteen nodes. Among these fourteen nodes, there are seven with the type annotation *assi*, two with type annotation *add*, and five with unique type annotations, *i.e.*, *start*, *goto*, *if*, *div* and *exit*. For each node, the d_{in} and d_{out} are calculated, too, to derive the complete NF. For instance, the node *start* (node 1) will be represented by the NF *start-0-1*, where *start* is NO_1 , 0 is d_{in} and 1 is d_{out} . Each unique NF is tallied to get the corresponding NF value. For *start-0-1*, the NF value is 1. Table II shows the set of NFs and their values extracted from Algorithm 1 using its CFG representation.

The feature based on path information (hereafter just PF) refers to the shortest routes from the start node to each node in the graph, as well as the shortest path from each node in the graph to the end node. This feature follows the form $NO_1 - NO_2 - NO_3 - \dots - NO_n$, where NO_n , as in the NF, denotes a specific operation statement in node n . The value of each PF is the number of occurrences of each path in the CFG. For instance, let us consider the labelled CFG of Algorithm 1. As Figure 2 right side shows, the path composed by the nodes the PF 1-2-3-4-5-6-7 and the nodes 1-2-3-4-5-6-11 can be denoted as *start-assi-assi-goto-assi-if-assi*. Therefore, its feature value is 2 since there are two paths represented by one type of PF. Table III shows the set of PFs and their associated PF values

extracted from Algorithm 1 using its CFG representation.

The second approach to extract features from CFG is by using *graph similarity measures*. Graph similarity refers to the process of determining the degree of similarity between two or more graphs. In particular, Kanewala *et al.* use Random Walk Kernel (RWK) and Graphlet Kernel (GK). RWK is the most-studied family of graph kernels [26]. It provides measure the similarity between two or more graphs based on the number of common walks in the graphs. The concept behind GK is to randomly sample tiny (connected) sub-graphs of size k , and using them to compare frequency distributions or to construct graph invariants.

TABLE I: Node operations (NO) in the control flow graph and corresponding labels (CL) for the annotation

| NO | CL | NO | CL | NO | CL | NO | CL |
|--------|-------|--------|--------|------|-------|------|------|
| + | add | − | sub | * | mul | / | div |
| , or | or | &, and | and | if | if | = | assi |
| == | eql | >= | geql | > | gt | <= | leql |
| < | lt | != | neql | := | start | % | rem |
| invoke | fcall | return | return | exit | exit | goto | goto |

Step 2.2 – MR labelling: The key idea of PMR is predicting whether a given method is suited for a particular MR by using binary classifiers. PMR uses supervising learning classification algorithms, *i.e.*, a labelled dataset is needed to provide examples for learning. Thus, after *Step 2.1 – Feature extraction*, the training dataset is created by manually labelling each method with applicable MRs. Depending on whether a specific MR does or does not satisfy the method, the method is labelled with 1 or 0 for this MR, respectively.

3) **Phase III – Training and testing:** This phase involves the use of one or more supervised machine learning (ML) algorithms, or a combination of them, to derive knowledge from the data. Three steps needs to be conducted.

Step 3.1 – Data split: This step is responsible for splitting the dataset into two subsets: a training set and a test set. The training set is used to create the prediction model, while the test set is used to evaluate the performance of the created prediction model.

Step 3.2 – Model creation refers to the process of building prediction models. Choosing a good modelling technique is vital for the training and prediction stage in any ML application, including the PMR approach. Kanewala *et al.* get the best results using the Support Vector Machine (SVM) technique.

Step 3.3 – Performance evaluation: This step measures the performance of the created prediction models. Performance

TABLE II: Node Features (NF) extracted from Algorithm 1 related to the nodes of its CFG representation

| NF | NF value | NF | NF value |
|-----------|----------|---------|----------|
| start-0-1 | 1 | if-2-2 | 1 |
| assi-1-1 | 7 | add-1-1 | 2 |
| goto-1-1 | 1 | div-1-1 | 1 |

TABLE III: Path Features (PF) extracted from Algorithm 1 related to the paths of its CFG representation

| PF | PF value |
|---|----------|
| Shortest path from the start node to each node | |
| start | 1 |
| start-assi | 1 |
| start-assi-assi | 1 |
| start-assi-assi-goto | 1 |
| start-assi-assi-goto-assi | 1 |
| start-assi-assi-goto-assi-if | 1 |
| start-assi-assi-goto-assi-if-assi | 2 |
| start-assi-assi-goto-assi-if-assi-assi | 2 |
| start-assi-assi-goto-assi-if-assi-assi-add | 1 |
| start-assi-assi-goto-assi-if-assi-assi-div | 1 |
| start-assi-assi-goto-assi-if-assi-assi-add-add | 1 |
| start-assi-assi-goto-assi-if-assi-assi-div-exit | 1 |
| Shortest path from each node to the end node | |
| assi-assi-goto-assi-if-assi-assi-div-exit | 1 |
| assi-goto-assi-if-assi-assi-div-exit | 1 |
| goto-assi-if-assi-assi-div-exit | 1 |
| assi-if-assi-assi-div-exit | 1 |
| if-assi-assi-div-exit | 1 |
| assi-assi-div-exit | 1 |
| assi-div-exit | 1 |
| div-exit | 1 |
| exit | 1 |
| assi-assi-add-add-if-assi-assi-div-exit | 1 |
| assi-add-add-if-assi-assi-div-exit | 1 |
| add-add-if-assi-assi-div-exit | 1 |
| add-if-assi-assi-div-exit | 1 |

measures are derived from the *Confusion Matrix*. Let A denote a classification output in which a specific MR_n satisfies the method m , and let A' denote a classification output in which a specific MR_n does not satisfy the method m , then A can be seen as the *positive class* and A' as the *negative class*. Using this notation, each standard performance measure is expressed as a function of the counts of elements in the *Confusion Matrix* defined as follows:

True Positive (TP): The actual MR of a method was A and the predicted was A . This represents a successful prediction.

True Negative (TN): The actual MR of a method was A' , the predicted was A' . This represents a successful prediction.

False Positive (FP): The actual MR was A' and the predicted was A . This represents an unsuccessful prediction.

False Negative (FN): The actual MR was A and the predicted was A' . This represents an unsuccessful prediction.

Accuracy is the ratio of successful predictions made to both classes and expressed as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

Precision (or positive predictive value) is the ratio of correct predictions made for class A and is shown in Equation (2):

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

Recall (or true positive rate, or sensitivity) is the ratio of successful predictions made to cases of class A

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

The f -measure statistic (or F1 score) is the harmonic mean of precision and recall:

$$f\text{-measure} = 2 \times \frac{Recall \times Precision}{Recall + Precision} \quad (4)$$

In addition to the aforementioned performance measures, the *Area Under Curve* (AUC) and the *Balanced Success Rate* (BSR) measures are also widely used. The AUC is the area under the curve that plots the False Positive Rate (FPR) against the True Positive Rate (TPR) at different points in $[0, 1]$. In binary classification problems, the BSR measure is calculated as the average of recall obtained on each class [27].

B. Metamorphic Relations

In the original study, Kanewala *et al.* use six MRs that had been suggested previously in other studies [7], [28]–[31]. Below we describe each MR in detail. The acronyms *stc* and *ftc* stand for *source test case* and *follow-up test case*, respectively. The input of the *stc* is an ordered set of non-negative numbers:

$$Input_{stc} = X_i, \dots, X_n \text{ where } X_i \geq 0, 0 \leq i \leq n$$

The outputs of the source and follow-up test cases are written as $Output_{stc}(X)$ and $Output_{ftc}(Y)$, respectively.

The MRs based on these inputs and outputs are:

MR₁: “Addition” (ADD). To get the *ftc* input, add a positive constant “ C ” to each element of the *stc* input, *i.e.*,

$$Input_{ftc} = X_1 + C, X_2 + C, X_3 + C, \dots, X_n + C,$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) \geq Output_{stc}(X)$$

MR₂: “Multiplication” (MUL). To get the *ftc* input, multiply each *stc* input element with a positive constant “ C ”, *i.e.*,

$$Input_{ftc} = X_1 * C, X_2 * C, \dots, X_n * C$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) \geq Output_{stc}(X)$$

MR₃: “Permutation” (PER). To get the *ftc* input, randomly permute the *stc* input elements, *e.g.*, like

$$Input_{ftc} = X_3, X_1, X_n, \dots, X_2$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) = Output_{stc}(X)$$

MR₄: “Inclusive” (INC). To get the *ftc* input, include a new element “ $X_{n+1} \geq 0$ ” to the *stc* input, *e.g.*, like

$$Input_{ftc} = X_1, X_2, X_3, \dots, X_n, X_{n+1}$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) \geq Output_{stc}(X)$$

MR₅: “Exclusive” (EXC). To get the *ftc* input, remove an element “ $X_{n-1} \geq 0$ ” from the *stc* input, *e.g.*, like

$$Input_{ftc} = X_1, X_2, X_3, \dots, X_{n-1}$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) \leq Output_{stc}(X)$$

MR₆: “Invertive” (INV). To get the *ftc* input, take the inverse of each *stc* input element $X_i > 0$, *i.e.*,

$$Input_{ftc} = 1/X_1, 1/X_2, 1/X_3, \dots, 1/X_n$$

Then the following output-relation must hold:

$$Output_{ftc}(Y) \leq Output_{stc}(X)$$

C. Dataset

In their original study, Kanewala *et al.* relied on a code corpus containing 100 Java methods that take numerical inputs and produce numerical outputs. The methods are from the open-source libraries *Colt Project* [32], which is an open-source library written for high-performance scientific and technical computing, *Apache Mahout* [33], which is a machine learning library, *Apache Commons Mathematics* [34], which is a Library of mathematics and statistics components, and *Java Collections* [35], which is a framework that provides an architecture to store and manipulate the group of objects. All of these libraries are written in Java.

To create a training dataset, Kanewala *et al.* manually labelled each method with the set of pre-defined MRs in a binary manner, *i.e.*, if MR_n matches a method m , then this method is assigned the label 1 for MR_n , otherwise it is 0.

Table IV reports the total number of methods that do and do not match a specific MR. One sees that more than half of the methods match with MRs denoted as ADD, MUL and INV, while approximately one third of the methods match with MRs denoted as PER, INC, and EXC.

TABLE IV: Total number of methods that match (✓) and do not match (✗) a specific MR

| MR | Change in the input | Output expected | ✓ | ✗ |
|-----|----------------------------------|-----------------------------|----|----|
| ADD | Add a positive constant | Increase or remain constant | 56 | 44 |
| MUL | Multiply by a positive constant | Increase or remain constant | 66 | 34 |
| PER | Permute the components | Remain constant | 33 | 67 |
| INC | Add a new element | Increase or remain constant | 34 | 66 |
| EXC | Remove an element | Decrease or remain constant | 32 | 68 |
| INV | Take the inverse of each element | Decrease or remain constant | 63 | 37 |

Table V reports how many methods have 0, 1, 2, ... 6 matching MRs and how those MRs are distributed in each case. 20 out of 100 methods have no matching MR, and only 9 out of 100 methods match with all six MRs simultaneously.

TABLE V: Total number of methods that have 0, 1, 2, ..., 6 matching MRs and their distributions

| No. MR [†] | No. Met [⊥] | ADD | MUL | PER | INC | EXC | INV |
|---------------------|----------------------|-----|-----|-----|-----|-----|-----|
| 0 | 20 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 8 | 2 | 3 | 0 | 2 | 0 | 1 |
| 2 | 7 | 3 | 4 | 2 | 1 | 1 | 3 |
| 3 | 23 | 19 | 17 | 5 | 5 | 5 | 18 |
| 4 | 26 | 16 | 26 | 16 | 10 | 10 | 26 |
| 5 | 7 | 7 | 7 | 1 | 7 | 7 | 6 |
| 6 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |

[†]Number of MRs that may apply to certain method, [⊥]Number of methods

TABLE VI: Labelled dataset [6]: symbol \checkmark denotes an MR-method match; symbol \times denotes that there is no match

| ID | Method Name | Library | Metamorphic Relation | | | | | | ID | Method Name | Library | Metamorphic Relation | | | | | |
|----|--------------------|---------|----------------------|--------------|--------------|--------------|--------------|--------------|-----|-------------------|---------|----------------------|--------------|--------------|--------------|--------------|--------------|
| | | | ADD | MUL | PER | INC | EXC | INV | | | | ADD | MUL | PER | INC | EXC | INV |
| 1 | add_values | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | 51 | find_median | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 2 | array_calc | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark | 52 | find_min | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 3 | array_copy | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark | 53 | g_Test | Math | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark |
| 4 | autoCorrelation | Colt | \times | \times | \times | \times | \times | \times | 54 | geometric_mean | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 5 | average | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark | 55 | get_array_value | Colle | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark |
| 6 | bi_SearchFromTo | Colt | \times | \times | \times | \checkmark | \times | \times | 56 | hamming_dist | Colle | \times | \times | \times | \checkmark | \checkmark | \checkmark |
| 7 | bubble | Math | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark | 57 | harmonicMean | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 8 | cal_AbsoluteDiff | Math | \checkmark | \checkmark | \times | \times | \times | \checkmark | 58 | insertion_sort | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 9 | cal_Diff | Math | \times | \times | \times | \times | \times | \times | 59 | kurtosis | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \times |
| 10 | chebyshevDist | Maho | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark | 60 | lag | Colt | \times | \checkmark | \times | \times | \times | \times |
| 11 | checkNonNegative | Math | \times | \checkmark | \checkmark | \checkmark | \times | \checkmark | 61 | manhattanDist | Maho | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark |
| 12 | checkPositive | Math | \times | \checkmark | \checkmark | \checkmark | \times | \checkmark | 62 | manhattanDist2 | Colle | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark |
| 13 | check_equal | Colle | \times | \times | \times | \times | \times | \times | 63 | max | Colt | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| 14 | check_eq_tolerance | Colle | \times | \times | \times | \times | \times | \times | 64 | meanDeviation | Colt | \checkmark | \times | \checkmark | \times | \times | \times |
| 15 | chiSquare | Math | \times | \checkmark | \times | \times | \times | \checkmark | 65 | mean_Diff | Math | \times | \times | \times | \times | \times | \times |
| 16 | clip | Colle | \times | \times | \times | \times | \times | \times | 66 | mean_abs_error | Colle | \times | \checkmark | \times | \times | \times | \checkmark |
| 17 | cnt_zeroes | Colle | \times | \times | \checkmark | \checkmark | \checkmark | \times | 67 | min | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 18 | canberraDist | Math | \times | \times | \times | \checkmark | \checkmark | \checkmark | 68 | partition | Math | \times | \times | \times | \times | \times | \times |
| 19 | cal_DividedDiff | Math | \checkmark | \times | \times | \times | \times | \times | 69 | polevl | Colt | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark |
| 20 | cosineDist | Maho | \times | \checkmark | \times | \times | \times | \times | 70 | pooledMean | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 21 | count_k | Colle | \times | \times | \checkmark | \checkmark | \checkmark | \times | 71 | pooledVariance | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 22 | count_non_zeroes | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | 72 | power | Colt | \times | \checkmark | \times | \times | \times | \checkmark |
| 23 | covariance | Colt | \checkmark | \times | \times | \times | \times | \times | 73 | product | Colt | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| 24 | dec | Maho | \times | \times | \times | \times | \times | \times | 74 | quantile | Colt | \checkmark | \checkmark | \times | \times | \checkmark | \checkmark |
| 25 | dec_array | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark | 75 | reverse | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 26 | Dist | Math | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark | 76 | safeNorm | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| 27 | DistInf | Math | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark | 77 | sampleKurtosis | Math | \checkmark | \times | \checkmark | \times | \times | \checkmark |
| 28 | dot_product | Colle | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark | 78 | sampleSkew | Colt | \checkmark | \times | \checkmark | \times | \times | \times |
| 29 | durbinWatson | Colt | \times | \checkmark | \times | \times | \times | \times | 79 | sampleVariance | Colt | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 30 | ebeAdd | Math | \checkmark | \checkmark | \times | \times | \times | \times | 80 | sampleWeightedVar | Colt | \times | \times | \times | \times | \times | \checkmark |
| 31 | ebeDivide | Math | \times | \times | \times | \times | \times | \times | 81 | scale | Colt | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 32 | ebeMultiply | Math | \checkmark | \checkmark | \times | \times | \times | \checkmark | 82 | s_add | Maho | \checkmark | \times | \times | \checkmark | \checkmark | \times |
| 33 | ebeSubtract | Math | \times | \times | \times | \times | \times | \times | 83 | selection_sort | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 34 | elemtWise_equal | Colle | \times | \times | \times | \times | \times | \times | 84 | sequential_search | Colle | \times | \times | \times | \checkmark | \checkmark | \times |
| 35 | elemtWise_max | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark | 85 | set_min_val | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 36 | elemtWise_min | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark | 86 | shell_sort | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 37 | elemtWise_not_eq | Colle | \times | \times | \times | \times | \times | \times | 87 | skew | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \times |
| 38 | entropy | Math | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \times | 88 | square | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 39 | equals | Math | \times | \times | \times | \times | \times | \times | 89 | standardize | Colle | \checkmark | \checkmark | \times | \times | \times | \times |
| 40 | errorRate | Maho | \times | \times | \times | \times | \times | \times | 90 | sum | Maho | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| 41 | euc_Dist | Math | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark | 91 | sumOfLogarithms | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark |
| 42 | evaluateHonors | Math | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark | 92 | sum_Power_Deviat | Colt | \times | \times | \times | \times | \times | \times |
| 43 | eval_Internal | Math | \times | \times | \times | \times | \times | \times | 93 | sum_labeled | Colt | \times | \times | \times | \times | \times | \times |
| 44 | evalNewton | Math | \times | \times | \times | \checkmark | \times | \times | 94 | tanimotoDist | Maho | \times | \times | \times | \times | \times | \times |
| 45 | evalWeightedProd | Math | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark | 95 | variance | Colle | \checkmark | \checkmark | \checkmark | \times | \times | \checkmark |
| 46 | find_diff | Colle | \times | \times | \times | \times | \times | \times | 96 | var_Difference | Colt | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 47 | find_euc_Dist | Colle | \times | \checkmark | \times | \checkmark | \checkmark | \checkmark | 97 | weightedMean | Colt | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 48 | find_magnitude | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | 98 | weightedRMS | Colt | \times | \times | \times | \times | \times | \times |
| 49 | find_max | Colle | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | \checkmark | 99 | weighted_average | Colle | \checkmark | \checkmark | \times | \times | \times | \checkmark |
| 50 | find_max2 | Colle | \checkmark | \checkmark | \times | \checkmark | \checkmark | \checkmark | 100 | winsorizedMean | Colt | \checkmark | \checkmark | \times | \times | \checkmark | \checkmark |

Colle: Java Collection, Maho: Apache Mahout, Math: Apache Commons Mathematics

Table VI shows the names of all methods used in the study and the library to which they belong. It also shows whether or not a specific MR matches the method. For better readability, we use the symbol \checkmark to denote that a specific MR matches, otherwise, we use the symbol \times . In total, 26 methods stem from the *Colt* project, 8 are from *Apache Mahout*, 25 are from *Apache Commons Mathematics*, and 41 methods are from *Java Collections*. These methods are provided by Kanewala *et al.*¹ in the form of CFG representation, using the graph description language format DOT, instead of source code.

D. Achieved Performance in Original Study

Kanewala *et al.* use features based on nodes and paths, as well as features based on graph similarity (RWK and GK), to

build 18 binary SVM models, *i.e.*, three models (each using different feature sets) per any of the six specific MRs. They use AUC and BSR to evaluate model performance and consider $AUC > 0.80$ to be a good classification performance. Among the eighteen trained SVM models, the most promising one was when RWK was used. The average performance for the six models using RWK was 0.87 in terms of AUC.

IV. REPLICATION METHODOLOGY

Our goal is to investigate whether the PMR approach of Kanewala *et al.* (i) can be replicated when using our own implementation of the pipeline for developing classifiers starting out from Java source code instead of CFG representations, (ii) classifiers trained on Java source code can be transferred to Python and C++ methods that have identical functionality,

¹<http://www.cs.colostate.edu/saxs/MRpred/functions.tar.gz>

and (iii) the PMR approach can be applied to Python and C++ code when classifiers are developed from scratch in the target programming languages. Each of these scenarios gives rise to a research question that we answer in our study. In all scenarios, we follow the PMR procedure, Figure 1, and we use the same set of six pre-defined MRs used in the original study, *i.e.*, Section III-B. However, we develop our own pipeline and create new datasets. For performance evaluation, we employ 10-fold stratified cross-validation, *i.e.*, the dataset is randomly partitioned into ten subgroups. The classifier is then built using nine subsets, with the 10th subset being used to evaluate the predictive model’s performance. This procedure is done ten times, with each of the ten subgroups being evaluated separately. The ten folds in stratified 10-fold cross-validation are partitioned in such a way that they include about the same proportion of classes as the original data set. The resulting overall performance is measured as the average of the ten cross-validation tests.

Our replication package containing results, scripts, models and datasets is available online².

A. Research Questions

We aim at answering the following research questions

- **RQ₁**: [Replicability] *How well do classifiers predict matching MRs for Java methods when using our processing and training pipeline starting from source code?*
- **RQ₂**: [Transferability] *How well do classifiers developed on Java code predict matching MRs for functionally equivalent methods implemented in Python and C++?*
- **RQ₃**: [Generalisability] *How well do classifiers predict matching methods for Python and C++ methods when developed from source code in the respective target languages?*

B. RQ₁: How well do classifiers predict matching MRs for Java methods when using our pipeline implementation?

In RQ1, we investigate the replicability of the PMR approach when starting out directly from Java source code (instead of CFG representations), performing all steps of feature extraction and re-generating the classifiers with a different ML package. In particular, we are interested in checking whether the classifiers developed by us achieve the same performance as published in [6]. Satisfactory results for RQ1 are the prerequisite for tackling RQ2 and RQ3.

To compare with the work of Kanewala *et al.*, we develop our own artefacts for *Phase II - Data preparation, Step 2.1 – Feature extraction*, and all the artefacts needed by *Phase III Training and testing*. Then we compare the results of SVM obtained by Kanewala *et al.* [6] who used *PyML Toolkit* [36] with our results achieved using Python *scikit-learn* library [37] with default parameter settings. For the comparison we use two datasets. The first dataset is the one used in the original study by Kanewala *et al.*, *i.e.*, the methods from Table VI. This dataset contains the CFG representations of the 100 Java

methods in DOT format. In the following, we call this dataset *DS_{JK}*. The second dataset contains the Java source code of the 100 methods from Table VI as contained in the open-source libraries. We call this dataset *DS_{JV}*.

When using *DS_{JK}* we apply the PMR approach from *Phase I - Step 1.2* through *Phase III - Step 3.3*, since this dataset already contains the CFG representation of each method in DOT format. When using *DS_{JV}*, we apply our entire pipeline implementing the PMR approach, *i.e.*, from *Phase I - Step 1.1* to *Phase III - Step 3.3*. To get the CFG representation in DOT format for the dataset *DS_{JV}*, we use the *soot*[25] Java framework, configured so that the output matches the CFGs of the original dataset. Then, we compare the performance of our classifiers against the results published by Kanewala *et al.* In particular, we compare BSR and AUC since these are the measures provided in the original study [6]. In addition, we provide the performance measures detailed in Section III-A3 - *Step 3.3 – Performance evaluation*.

C. RQ₂: How well do classifiers developed on Java code predict matching MRs for functionally equivalent methods implemented in Python and C++?

In RQ2, we check whether classifiers developed with the PMR approach from Java methods achieve the same performance when applied to methods with identical functionality but implemented in Python or C++. We chose Python and C++ because both are popular and widely used programming languages supporting a broad range of applications [38].

For this experiment, we created two new datasets containing source code of methods written in Python and in C++. The methods in each dataset are functionally identical to that of the 100 Java methods described in Table VI. The corresponding method implementations were either retrieved from the *NumPy* package for scientific computing in case of Python, the *Blinz++* high-performance library for scientific computing in case of C++, or they were implemented in Python/C++ by the authors if not present in these libraries. Because the functionality of the Python and C++ methods is equivalent to the functionality of the Java methods, we can assume that exactly the same MRs that match the Java methods match the corresponding Python and C++ methods. The dataset named *DS_{PY}* contains the Python methods, and the dataset named *DS_{C++}* contains the C++ methods. To get the graph representation in DOT format, Section III-A1: *Phase I - Step CFG generation*, for the methods written in Python, we use the Python package *pycfg* [39], and for the methods written in C++, we use *Goblint*[40], [41].

D. RQ₃: How well do classifiers predict matching methods for Python and C++ methods when developed from source code in the respective target languages?

Finally, in RQ3, we check whether the PMR approach works for Python and C++ code similarly well as it does for Java code, if we develop the classifiers for each target language from scratch. Thus, we train SVM models using each dataset,

²<https://github.com/aduquet/RENE-PredictingMetamorphicRelations>

i.e., DS_{PY} and DS_{C++} . We compare the performance of the new classifiers against the results obtained in RQ1 and RQ2.

V. RESULTS AND DISCUSSION

A. RQ₁ How well do classifiers predict matching MRs for Java methods when developed from source code using our pipeline?

Table VII shows the performance of our PMR implementation for both Java datasets, DS_{JK} and DS_{JV} . Overall, regardless of the feature extraction technique used, the results are fairly close. This can be seen in the Error column, which displays the difference in performance between DS_{JK} and DS_{JV} for each MR. The most negative value is -0.104 (Accuracy of ADD) while the farthest positive value is 0.148 (BSR of INV). This indicates that the classifiers developed by us are consistent for Java code independent from the starting point of the model development (CFG vs. source code).

Table VIII shows how the performance of our PMR implementation compares to the performance obtained by Kanewala *et al.* in terms of AUC and BSR. As can be seen from the Error column, our results are close to those obtained in the original study. The Error range is $[-0.093, 0.061]$ for AUC and $[-0.118, 0.117]$ for BSR. From combining the results shown in Table VIII with those shown in Table VII we conclude that our implementation of PMR achieves similar performance as reported in [6] even when starting out from source code.

With regards to *replicability* (RQ1), our results indicate that we can achieve similar results as Kanewala *et al.* when re-implementing the PMR approach no matter whether we start the modelling process from source code or from CFG representations.

B. RQ₂ How well do classifiers developed on Java code predict matching MRs for functionally equivalent methods implemented in Python and C++?

Table IX reports on the performance when using classifiers, developed starting out from the DS_{JV} dataset, to predict matching MRs for methods contained in the DS_{PY} and DS_{C++} datasets. The assumption behind applying a classifier built on Java code to methods that are functionally equivalent but implemented in a different programming language is that the CFG representations from which the features in the SVM models are taken would be similar enough to achieve similar classification performance as when applied to Java methods.

However, as shown in Table IX, the performance is low for all performance measures and for both Python and C++. No measure is greater than 0.689 . This result suggests that the representation of the CFGs of the Python and C++ methods to which the feature extraction algorithm is applied are more different from the CFGs of the Java methods than expected. This can be explained due to the language-specific CFG generators that we used as well as differences in the way how the methods (with identical functionality) are implemented in different programming languages.

With regards to *transferability* (RQ2), our results suggest that classifiers trained on a dataset containing methods in one programming language (Java) have reduced performance when applied to datasets with functionally equivalent methods implemented in a different programming language (Python, C++). Hence, classification models built according to the proposed PMR approach may not be transferable across languages.

C. RQ₃ How well do classifiers predict matching methods for Python and C++ methods when developed from source code in the respective target languages?

Table X reports the results of using the PMR approach to develop classifiers separately for each programming language (Python and C++). Comparing Table IX and Table X indicates that the performance improves remarkably when using models that are trained specifically to also consider the implementation characteristics stemming from the different programming languages. Even though the performance has improved by developing language specific classifiers, the results for Python and C++ are generally below the results achieved for Java, with the results for C++ being consistently the worst.

With regards to *generalizability* (RQ3), our results suggest the PMR approach can be applied for different programming languages when the classifiers are re-trained on the specific target language. The slightly lower performance, esp. for C++, needs further exploration of the data and the choice of model parameter settings (tuning).

D. Threats to Validity

In the context of our study, two types of threats to validity are most relevant: threats to internal and external validity.

To achieve internal validity, we used the same set of methods and of MRs as in Kanewala *et al.* [6]. For the Python and C++ datasets, we carefully checked functional equivalence of the methods with those in the original Java dataset. Given functional equivalence of the methods, we assume that the matching MRs are identical for each of the three chosen programming languages. However, this has not been verified. It is unlikely but possible that some methods have slight differences in the set of matching MRs due to the programming language. Another potential validity threat in our study is that we recreated all steps of the PMR approach using different machine learning libraries with potentially different parameter settings. However, the performance measures in RQ1 (Table VII) align well with the results reported in the original study. This suggests that we have understood how to correctly build the classifiers in our replication.

Regarding external validity, our study uses the same methods as in the original study but implemented in different

TABLE VII: PMR performance achieved by our classifiers when starting from DS_{JV} and DS_{JK}

| MR | Feat [⊥] | Performance measurements | | | | | | | | | | | | | | | | | |
|-----|-------------------|--------------------------|--------------|--------------------|--------------|--------------|--------------------|--------------|--------------|--------------------|--------------|--------------|--------------------|--------------|--------------|--------------------|--------------|--------------|--------------------|
| | | Accuracy | | | Precision | | | Recall | | | f-measure | | | AUC | | | BSR | | |
| | | DS_{JK} | DS_{JV} | Error [±] | DS_{JK} | DS_{JV} | Error [±] | DS_{JK} | DS_{JV} | Error [±] | DS_{JK} | DS_{JV} | Error [±] | DS_{JK} | DS_{JV} | Error [±] | DS_{JK} | DS_{JV} | Error [±] |
| ADD | NF-NP | 0.802 | 0.787 | 0.015 | 0.786 | 0.751 | 0.035 | 0.812 | 0.704 | 0.108 | 0.773 | 0.775 | -0.002 | 0.837 | 0.827 | 0.010 | 0.768 | 0.785 | -0.017 |
| | GK | 0.712 | 0.816 | -0.104 | 0.702 | 0.732 | -0.030 | 0.717 | 0.758 | -0.041 | 0.744 | 0.712 | 0.032 | 0.769 | 0.707 | 0.062 | 0.737 | 0.729 | 0.008 |
| | RWK | 0.851 | 0.86 | -0.009 | 0.836 | 0.712 | 0.124 | 0.771 | 0.791 | -0.020 | 0.786 | 0.785 | 0.001 | 0.905 | 0.877 | 0.028 | 0.843 | 0.829 | 0.014 |
| MUL | NF-NP | 0.712 | 0.688 | 0.024 | 0.672 | 0.689 | -0.017 | 0.685 | 0.661 | 0.024 | 0.657 | 0.705 | -0.048 | 0.742 | 0.734 | 0.008 | 0.631 | 0.654 | -0.023 |
| | GK | 0.663 | 0.641 | 0.022 | 0.714 | 0.732 | -0.018 | 0.697 | 0.758 | -0.061 | 0.676 | 0.733 | -0.057 | 0.775 | 0.730 | 0.045 | 0.689 | 0.657 | 0.032 |
| | RWK | 0.789 | 0.695 | 0.094 | 0.666 | 0.706 | -0.040 | 0.693 | 0.797 | -0.104 | 0.660 | 0.676 | -0.016 | 0.846 | 0.820 | 0.026 | 0.774 | 0.739 | 0.035 |
| PER | NF-NP | 0.838 | 0.840 | -0.002 | 0.860 | 0.883 | -0.023 | 0.835 | 0.846 | -0.011 | 0.855 | 0.790 | 0.065 | 0.945 | 0.925 | 0.020 | 0.847 | 0.813 | 0.034 |
| | GK | 0.834 | 0.826 | 0.008 | 0.888 | 0.819 | 0.069 | 0.823 | 0.845 | -0.022 | 0.864 | 0.79 | 0.074 | 0.872 | 0.811 | 0.061 | 0.853 | 0.839 | 0.014 |
| | RWK | 0.916 | 0.918 | -0.002 | 0.917 | 0.827 | 0.090 | 0.878 | 0.835 | 0.043 | 0.877 | 0.893 | -0.016 | 0.963 | 0.944 | 0.019 | 0.757 | 0.793 | -0.036 |
| INC | NF-NP | 0.792 | 0.807 | -0.015 | 0.847 | 0.822 | 0.025 | 0.837 | 0.837 | 0.000 | 0.776 | 0.759 | 0.017 | 0.845 | 0.852 | -0.007 | 0.793 | 0.786 | 0.007 |
| | GK | 0.752 | 0.788 | -0.036 | 0.721 | 0.781 | -0.060 | 0.790 | 0.776 | 0.014 | 0.783 | 0.718 | 0.065 | 0.850 | 0.882 | -0.032 | 0.762 | 0.744 | 0.018 |
| | RWK | 0.799 | 0.839 | -0.040 | 0.832 | 0.792 | 0.040 | 0.800 | 0.773 | 0.027 | 0.854 | 0.764 | 0.090 | 0.862 | 0.821 | 0.041 | 0.673 | 0.654 | 0.019 |
| EXC | NF-NP | 0.763 | 0.753 | 0.010 | 0.772 | 0.783 | -0.011 | 0.778 | 0.759 | 0.019 | 0.762 | 0.789 | -0.027 | 0.768 | 0.755 | 0.013 | 0.868 | 0.839 | 0.029 |
| | GK | 0.816 | 0.787 | 0.029 | 0.816 | 0.861 | -0.045 | 0.849 | 0.890 | -0.041 | 0.871 | 0.790 | 0.081 | 0.873 | 0.840 | 0.003 | 0.758 | 0.755 | 0.003 |
| | RWK | 0.774 | 0.725 | 0.049 | 0.757 | 0.743 | 0.014 | 0.757 | 0.741 | 0.016 | 0.769 | 0.744 | 0.025 | 0.731 | 0.727 | 0.004 | 0.79 | 0.757 | 0.033 |
| INV | NF-NP | 0.714 | 0.705 | 0.009 | 0.674 | 0.659 | 0.015 | 0.702 | 0.671 | 0.031 | 0.675 | 0.694 | -0.019 | 0.905 | 0.917 | -0.012 | 0.656 | 0.661 | -0.005 |
| | GK | 0.778 | 0.759 | 0.019 | 0.765 | 0.769 | -0.004 | 0.738 | 0.737 | 0.001 | 0.760 | 0.721 | 0.039 | 0.671 | 0.670 | 0.001 | 0.679 | 0.655 | 0.024 |
| | RWK | 0.651 | 0.585 | 0.066 | 0.610 | 0.659 | -0.049 | 0.643 | 0.639 | 0.004 | 0.675 | 0.653 | 0.022 | 0.760 | 0.766 | -0.006 | 0.787 | 0.639 | 0.148 |

[⊥]Feature extraction approach, [±]Error ($DS_{JK} - DS_{JV}$), **NF-PF**: Node Feature - Path Feature, **GK**: Graphnet Kernel, **RWK**: Random Walk Kernel

 TABLE VIII: Comparison of PMR performance (AUC and BSR) achieved by Kanewala *et al.* and when using classifiers developed by us starting from DS_{JK}

| MR | Feat [⊥] | Performance measurements | | | | | |
|-----|-------------------|--------------------------|--------------|--------------------|-------------|--------------|--------------------|
| | | [6] | DS_{JK} | Error [±] | [6] | DS_{JK} | Error [±] |
| ADD | NF-PF | 0.81 | 0.837 | -0.027 | 0.77 | 0.768 | 0.002 |
| | GK | 0.83 | 0.769 | 0.061 | 0.79 | 0.737 | 0.053 |
| | RWK | 0.92 | 0.905 | 0.015 | 0.85 | 0.843 | 0.007 |
| MUL | NF-PF | 0.73 | 0.742 | -0.012 | 0.65 | 0.631 | 0.019 |
| | GK | 0.78 | 0.775 | 0.005 | 0.69 | 0.689 | 0.001 |
| | RWK | 0.83 | 0.846 | -0.016 | 0.74 | 0.774 | -0.034 |
| PER | NF-PF | 0.93 | 0.945 | -0.015 | 0.83 | 0.847 | -0.017 |
| | GK | 0.91 | 0.872 | 0.038 | 0.83 | 0.853 | -0.023 |
| | RWK | 0.95 | 0.963 | -0.013 | 0.87 | 0.757 | 0.113 |
| INC | NF-PF | 0.84 | 0.845 | -0.005 | 0.80 | 0.793 | 0.007 |
| | GK | 0.88 | 0.850 | 0.030 | 0.75 | 0.762 | -0.012 |
| | RWK | 0.89 | 0.862 | 0.028 | 0.79 | 0.673 | 0.117 |
| EXC | NF-PF | 0.78 | 0.768 | 0.012 | 0.75 | 0.868 | -0.118 |
| | GK | 0.78 | 0.873 | -0.093 | 0.74 | 0.758 | -0.018 |
| | RWK | 0.90 | 0.731 | 0.169 | 0.79 | 0.790 | 0.000 |
| INV | NF-PF | 0.84 | 0.905 | -0.065 | 0.64 | 0.656 | -0.016 |
| | GK | 0.68 | 0.671 | 0.009 | 0.66 | 0.679 | -0.019 |
| | RWK | 0.76 | 0.769 | -0.009 | 0.74 | 0.787 | -0.047 |

[⊥]Feature extraction approach, [±]Error ($[6] - DS_{JK}$), [6]: Kanewala *et al.*

NF-PF: Node and Path Feature, **GK**: Graphnet Kernel

RWK: Random Walk Kernel

programming languages. For the sake of generalisability, it would have been preferable to include additional methods to overcome any potential bias introduced by the selection of methods in the original study. As a consequence, our replication cannot determine the actual scope of the effectiveness of the PMR approach.

E. Remarks on General Relevance

When assembling the Python and C++ datasets containing functionally equivalent methods for our replication, we

identified the issue that such methods tend to be rare and are usually only found in specific domains such as libraries for mathematical computations. In the original study, a fully labelled dataset was used containing a high number of methods (80%) with matching MRs. Only 20% of the methods are not related to any of the supported MRs. How realistic is this distribution? Since the pre-defined set of MRs is rather small and only applies to methods with a very specific signature (mainly functions that take numerical inputs and produce numerical outputs), it is unlikely that one will find an equal share of such methods in real-world applications. In particular since such methods are often already provided as part of existing, dedicated libraries (e.g., Apache Commons or NumPy). If, as we assume, the share of matching methods in newly developed real-world application is very small and given that the effort for developing language-specific classifiers is comparably high, the practical relevance of the proposed approach seems to be limited.

Furthermore, the proposed PMR approach uses features extracted from individual methods and it is therefore tied to the level of unit testing. A generalisation of the approach beyond unit testing, e.g., by transferring it to system level testing does not seem possible.

VI. CONCLUSION

We closely as well as conceptually replicated the study of Kanewala *et al.* [6]. First, we reproduced the PMR approach using our own implementation of the pipeline for feature extraction and training classifiers by starting out from Java source code and creating corresponding CFGs. We showed that our classifiers perform equally well as in the original study indicating a successful replication as basis for further experiments. Second, we checked transferability of classifiers trained on methods implemented in Java to other programming languages (Python and C++). We found that the performance decreases too much to consider this approach feasible. This is caused by programming language-specific implementation

TABLE IX: Performance of SVM models when trained with DS_{JV} and tested with DS_{PY} and DS_{C++}

| MR | Feat [⊥] | Performance measurements | | | | | | | | | | | |
|-----|-------------------|--------------------------|------------|-----------|------------|-----------|------------|-----------|------------|-----------|------------|-----------|------------|
| | | Accuracy | | Precision | | Recall | | f-measure | | AUC | | BSR | |
| | | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} |
| ADD | NF-PF | 0.575 | 0.459 | 0.572 | 0.522 | 0.555 | 0.551 | 0.554 | 0.473 | 0.551 | 0.529 | 0.563 | 0.466 |
| | GK | 0.564 | 0.447 | 0.532 | 0.426 | 0.543 | 0.470 | 0.531 | 0.473 | 0.547 | 0.452 | 0.561 | 0.427 |
| | RWK | 0.544 | 0.468 | 0.526 | 0.474 | 0.593 | 0.403 | 0.500 | 0.483 | 0.550 | 0.466 | 0.503 | 0.414 |
| MUL | NF-PF | 0.494 | 0.588 | 0.627 | 0.596 | 0.495 | 0.639 | 0.522 | 0.658 | 0.623 | 0.574 | 0.652 | 0.648 |
| | GK | 0.460 | 0.472 | 0.463 | 0.436 | 0.477 | 0.392 | 0.479 | 0.400 | 0.480 | 0.431 | 0.475 | 0.478 |
| | RWK | 0.499 | 0.388 | 0.499 | 0.388 | 0.492 | 0.387 | 0.488 | 0.393 | 0.490 | 0.393 | 0.499 | 0.384 |
| PER | NF-PF | 0.521 | 0.445 | 0.503 | 0.403 | 0.517 | 0.411 | 0.507 | 0.570 | 0.535 | 0.519 | 0.542 | 0.545 |
| | GK | 0.520 | 0.458 | 0.525 | 0.398 | 0.563 | 0.431 | 0.546 | 0.392 | 0.499 | 0.399 | 0.535 | 0.454 |
| | RWK | 0.515 | 0.424 | 0.515 | 0.540 | 0.503 | 0.471 | 0.517 | 0.496 | 0.516 | 0.413 | 0.532 | 0.514 |
| INC | NF-PF | 0.579 | 0.578 | 0.576 | 0.527 | 0.580 | 0.496 | 0.580 | 0.590 | 0.597 | 0.603 | 0.577 | 0.477 |
| | GK | 0.578 | 0.518 | 0.588 | 0.501 | 0.597 | 0.576 | 0.579 | 0.476 | 0.582 | 0.558 | 0.594 | 0.587 |
| | RWK | 0.536 | 0.521 | 0.525 | 0.444 | 0.508 | 0.528 | 0.536 | 0.526 | 0.512 | 0.478 | 0.500 | 0.486 |
| EXC | NF-PF | 0.637 | 0.440 | 0.639 | 0.497 | 0.535 | 0.502 | 0.591 | 0.507 | 0.600 | 0.525 | 0.639 | 0.478 |
| | GK | 0.597 | 0.561 | 0.648 | 0.506 | 0.592 | 0.552 | 0.649 | 0.494 | 0.610 | 0.467 | 0.658 | 0.492 |
| | RWK | 0.568 | 0.548 | 0.579 | 0.564 | 0.613 | 0.610 | 0.579 | 0.627 | 0.570 | 0.590 | 0.637 | 0.562 |
| INV | NF-PF | 0.531 | 0.493 | 0.534 | 0.422 | 0.525 | 0.433 | 0.502 | 0.471 | 0.514 | 0.411 | 0.512 | 0.491 |
| | GK | 0.472 | 0.411 | 0.478 | 0.395 | 0.473 | 0.401 | 0.477 | 0.421 | 0.468 | 0.403 | 0.469 | 0.459 |
| | RWK | 0.470 | 0.417 | 0.507 | 0.378 | 0.529 | 0.400 | 0.465 | 0.402 | 0.435 | 0.333 | 0.461 | 0.352 |

[⊥]Feature extraction approach, **NF-PF**: Node Feature - Path Feature, **GK**: Graphnet Kernel, **RWK**: Random Walk Kernel

TABLE X: Performance of SVM models for DS_{PY} and DS_{C++} datasets

| MR | Feat [⊥] | Performance measurements | | | | | | | | | | | |
|-----|-------------------|--------------------------|------------|-----------|------------|-----------|------------|-----------|------------|-----------|------------|-----------|------------|
| | | Accuracy | | Precision | | Recall | | f-measure | | AUC | | BSR | |
| | | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} | DS_{PY} | DS_{C++} |
| ADD | NF-PF | 0.706 | 0.577 | 0.742 | 0.660 | 0.748 | 0.590 | 0.683 | 0.645 | 0.723 | 0.691 | 0.760 | 0.653 |
| | GK | 0.724 | 0.599 | 0.671 | 0.708 | 0.713 | 0.655 | 0.757 | 0.606 | 0.730 | 0.652 | 0.798 | 0.667 |
| | RWK | 0.737 | 0.738 | 0.653 | 0.720 | 0.699 | 0.797 | 0.668 | 0.730 | 0.693 | 0.750 | 0.726 | 0.725 |
| MUL | NF-PF | 0.670 | 0.611 | 0.652 | 0.623 | 0.701 | 0.584 | 0.787 | 0.688 | 0.746 | 0.594 | 0.656 | 0.580 |
| | GK | 0.613 | 0.658 | 0.627 | 0.657 | 0.659 | 0.648 | 0.643 | 0.561 | 0.663 | 0.607 | 0.663 | 0.625 |
| | RWK | 0.727 | 0.795 | 0.732 | 0.742 | 0.640 | 0.685 | 0.726 | 0.715 | 0.686 | 0.735 | 0.721 | 0.677 |
| PER | NF-PF | 0.818 | 0.732 | 0.822 | 0.702 | 0.820 | 0.778 | 0.755 | 0.763 | 0.769 | 0.754 | 0.865 | 0.754 |
| | GK | 0.835 | 0.777 | 0.785 | 0.725 | 0.820 | 0.830 | 0.802 | 0.752 | 0.797 | 0.717 | 0.829 | 0.694 |
| | RWK | 0.869 | 0.824 | 0.856 | 0.853 | 0.811 | 0.877 | 0.862 | 0.755 | 0.796 | 0.735 | 0.798 | 0.840 |
| INC | NF-PF | 0.746 | 0.677 | 0.734 | 0.684 | 0.789 | 0.661 | 0.796 | 0.705 | 0.789 | 0.647 | 0.792 | 0.619 |
| | GK | 0.671 | 0.754 | 0.659 | 0.713 | 0.685 | 0.756 | 0.685 | 0.781 | 0.682 | 0.772 | 0.681 | 0.672 |
| | RWK | 0.785 | 0.760 | 0.793 | 0.721 | 0.808 | 0.784 | 0.746 | 0.682 | 0.804 | 0.753 | 0.793 | 0.694 |
| EXC | NF-PF | 0.734 | 0.635 | 0.694 | 0.649 | 0.713 | 0.652 | 0.725 | 0.690 | 0.715 | 0.732 | 0.737 | 0.732 |
| | GK | 0.752 | 0.698 | 0.735 | 0.681 | 0.703 | 0.745 | 0.725 | 0.742 | 0.764 | 0.760 | 0.734 | 0.673 |
| | RWK | 0.791 | 0.805 | 0.805 | 0.834 | 0.782 | 0.794 | 0.788 | 0.786 | 0.808 | 0.811 | 0.780 | 0.753 |
| INV | NF-PF | 0.606 | 0.571 | 0.671 | 0.543 | 0.668 | 0.539 | 0.661 | 0.594 | 0.673 | 0.559 | 0.672 | 0.538 |
| | GK | 0.582 | 0.595 | 0.598 | 0.620 | 0.589 | 0.586 | 0.558 | 0.579 | 0.604 | 0.633 | 0.599 | 0.624 |
| | RWK | 0.683 | 0.611 | 0.673 | 0.629 | 0.717 | 0.614 | 0.648 | 0.649 | 0.675 | 0.711 | 0.712 | 0.708 |

[⊥]Feature extraction approach, **NF-PF**: Node Feature - Path Feature, **GK**: Graphnet Kernel, **RWK**: Random Walk Kernel

details, despite relying only on features extracted from the abstract CFG representation of the methods. Third, we demonstrated that the PMR approach can be generalized. When re-training the classifiers from scratch on Python and C++ source code, the performance we achieved was almost comparable to those from classifiers trained on Java code.

All artefacts created by us as well as all results are available in a replication package.

ACKNOWLEDGEMENT

This research was partly funded by the Estonian Center of Excellence in ICT research (EXCITE), the European Regional Development Fund, the IT Academy Programme for ICT Research Development, the Austrian ministries BMVIT and BMDW, the State of Upper Austria under the COMET (Competence Centers for Excellent Technologies) program managed by FFG, and grant PRG1226 of the Estonian Research Council.

REFERENCES

- [1] T. Y. Chen, S. C. Cheung, and S. M. Yiu, "Metamorphic testing: A new approach for generating next test cases," *Department of Computer Science, Hong Kong University of Science and Technology, Hong Kong, Tech. Rep. HKUST-CS98-01*, 1998.
- [2] A. Duque-Torres, A. Shalygina, D. Pfahl, and R. Ramler, "Using rule mining for automatic test oracle generation," in *8th International Workshop on Quantitative Approaches to Software Quality (QuASoQ)*, ser. QuASoQ'20, 2020.
- [3] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The oracle problem in software testing: A survey," *IEEE Transactions on Software Engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [4] H. Liu, F.-C. Kuo, D. Towey, and T. Y. Chen, "How effectively does metamorphic testing alleviate the oracle problem?" *IEEE Transactions on Software Engineering*, vol. 40, no. 1, pp. 4–22, 2014.
- [5] Z. Q. Zhou, L. Sun, T. Y. Chen, and D. Towey, "Metamorphic relations for enhancing system understanding and use," *IEEE Transactions on Software Engineering*, vol. 46, no. 10, pp. 1120–1154, 2020.
- [6] U. Kanewala, J. M. Bieman, and A. Ben-Hur, "Predicting metamorphic relations for testing scientific software: A machine learning approach using graph kernels," *Software testing, verification and reliability*, vol. 26, no. 3, pp. 245–269, 2016.
- [7] U. Kanewala and J. M. Bieman, "Using machine learning techniques to detect metamorphic relations for programs without test oracles," in *IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)*, 2013, pp. 1–10.
- [8] M. Shepperd, N. Ajiénka, and S. Counsell, "The role and value of replication in empirical software engineering results," *Information and Software Technology*, vol. 99, pp. 120–132, 2018.
- [9] O. S. Gómez, N. Juristo, and S. Vegas, "Understanding replication of experiments in software engineering: A classification," *Information and Software Technology*, vol. 56, no. 8, pp. 1033–1048, 2014.
- [10] J. C. Carver, "Towards reporting guidelines for experimental replications: A proposal," in *1st international workshop on replication in empirical software engineering*, Citeseer, vol. 1, 2010, pp. 1–4.
- [11] M. LLC. (2020). "Artifact review and badging - current," [Online]. Available: <https://www.acm.org/publications/policies/artifact-review-and-badging-current> (visited on 09/21/2021).
- [12] M. Zhang, Y. Zhang, L. Zhang, C. Liu, and S. Khurshid, "Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems," in *33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, IEEE, 2018, pp. 132–142.
- [13] Z. Q. Zhou and L. Sun, "Metamorphic testing of driverless cars," *Communications of the ACM*, vol. 62, no. 3, pp. 61–67, Feb. 2019.
- [14] P. C. Canizares, A. Núñez, J. de Lara, and L. Llana, "MT-EA4Cloud: A methodology for testing and optimising energy-aware cloud systems," *Journal of Systems and Software*, vol. 163, p. 110522, 2020.
- [15] Z. Zhang, D. Towey, Z. Ying, Y. Zhang, and Z. Q. Zhou, "MT4NS: Metamorphic testing for network scanning," in *6th IEEE/ACM International Workshop on Metamorphic Testing (MET)*, ser. MET'21, 2021, pp. 17–23.
- [16] M. Srinivasan, M. P. Shahri, I. Kahanda, and U. Kanewala, "Quality assurance of bioinformatics software: A case study of testing a biomedical text processing tool using metamorphic testing," in *IEEE/ACM 3rd International Workshop on Metamorphic Testing (MET)*, ser. MET'18, Gothenburg, Sweden: Association for Computing Machinery, 2018, pp. 26–33.
- [17] M. P. Shahri, M. Srinivasan, G. Reynolds, D. Bimczok, I. Kahanda, and U. Kanewala, "Metamorphic testing for quality assurance of protein function prediction tools," in *IEEE International Conference On Artificial Intelligence Testing (AITest)*, IEEE, 2019, pp. 140–148.
- [18] Z. Peng, U. Kanewala, and N. Niu, "Contextual understanding and improvement of metamorphic testing in scientific software development," in *15th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, 2021, pp. 1–6.
- [19] X. Lin, M. Simon, and N. Niu, "Exploratory metamorphic testing for scientific software," *Computing in Science Engineering*, vol. 22, no. 2, pp. 78–87, 2020.
- [20] T. Y. Chen, D. Huang, T. Tse, and Z. Q. Zhou, "Case studies on the selection of useful relations in metamorphic testing," in *4th Ibero-American Symposium on Software Engineering and Knowledge Engineering (JIISIC 2004)*, Citeseer, 2004, pp. 569–583.
- [21] H. Liu, X. Liu, and T. Y. Chen, "A new method for constructing metamorphic relations," in *12th International Conference on Quality Software*, 2012, pp. 59–68.
- [22] J. Zhang, J. Chen, D. Hao, Y. Xiong, B. Xie, L. Zhang, and H. Mei, "Search-based inference of polynomial metamorphic relations," in *29th ACM/IEEE International Conference on Automated Software Engineering*, ser. ASE'14, Vasteras, Sweden, 2014, pp. 701–712.
- [23] T. Y. Chen, P.-L. Poon, and X. Xie, "METRIC: METAmorphic Relation Identification based on the Category-choice framework," *Journal of Systems and Software*, vol. 116, pp. 177–190, 2016.
- [24] C.-A. Sun, A. Fu, P.-L. Poon, X. Xie, H. Liu, and T. Y. Chen, "METRIC+: A metamorphic relation identification technique based on input plus output domains," *IEEE Transactions on Software Engineering*, vol. 47, no. 9, pp. 1764–1785, 2021.
- [25] R. Vallée-Rai, P. Co, E. Gagnon, L. Hendren, P. Lam, and V. Sundaresan, "Soot: A java bytecode optimization framework," in *CASCON First Decade High Impact Papers*, ser. CASCON '10, Toronto, Ontario, Canada: IBM Corp., 2010, pp. 214–224.
- [26] T. Gärtner, P. Flach, and S. Wrobel, "On graph kernels: Hardness results and efficient alternatives," in *Learning Theory and Kernel Machines*, Springer Berlin Heidelberg, 2003, pp. 129–143.
- [27] A. Ben-Hur and J. Weston, "A user's guide to support vector machines," in *Data Mining Techniques for the Life Sciences*. Humana Press, 2010, pp. 223–239.
- [28] U. Kanewala, "Techniques for automatic detection of metamorphic relations," in *IEEE 7th International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, 2014, pp. 237–238.
- [29] B. Hardin and U. Kanewala, "Using semi-supervised learning for predicting metamorphic relations," in *3rd IEEE/ACM International Workshop on Metamorphic Testing (MET)*, ser. MET'18, 2018, pp. 14–17.
- [30] K. Rahman and U. Kanewala, "Predicting metamorphic relations for matrix calculation programs," in *3rd IEEE/ACM International Workshop on Metamorphic Testing (MET)*, ser. MET'18, 2018, pp. 10–13.
- [31] K. Rahman, I. Kahanda, and U. Kanewala, "MRpredT: Using text mining for metamorphic relation prediction," in *42nd IEEE/ACM International Conference on Software Engineering Workshops (ICSEW)*, 2020, pp. 420–424.
- [32] *Colt project*, <http://acs.lbl.gov/software/colt/>, Accessed: 2021-09-21.
- [33] *Apache mahout*, <https://mahout.apache.org/>, Accessed: 2021-09-21.
- [34] *Apache commons mathematic*, <http://commons.apache.org/proper/commons-math/>, Accessed: 2021-09-21.
- [35] *Java collections*, <https://docs.oracle.com/javase/8/docs/technotes/guides/collections/overview.html>, Accessed: 2021-09-21.
- [36] *Pyml toolkit*, <http://pyml.sourceforge.net/>, Accessed: 2021-09-21.
- [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine Learning in Python," in *Journal of Machine Learning Research*, vol. 12, 2011, pp. 2825–2830.
- [38] S. Cass, "The top programming languages: Our latest rankings put python on top-again-[careers]," *IEEE Spectrum*, vol. 57, no. 8, pp. 22–22, 2020.
- [39] *Pycfg*, <https://pypi.org/project/pycfg/>, Accessed: 2021-09-21.
- [40] V. Vojdani, K. Apinis, V. Rõtov, H. Seidl, V. Vene, and R. Vogler, "Static race detection for device drivers: The goblin approach," in *31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, ser. ASE'16, ACM, 2016, pp. 391–402.
- [41] *Goblint github*, <https://github.com/goblint/analyzer>, Accessed: 2021-09-21.