

Semantic Clone Detection via Probabilistic Software Modeling

Hannes Thaller ^{*}, Lukas Linsbauer [†], Brent van Bladel [‡], and Alexander Egyed ^{*}

^{*} Johannes Kepler University Linz, Austria

{hannes.thaller, alexander.egyed}@jku.at

[†] Technical University of Braunschweig, Germany

l.linsbauer@tu-braunschweig.de

[‡] University of Antwerp, Belgium

brent.vandbladel@uantwerpen.be

Abstract—Semantic clone detection is the process of finding program elements with similar or equal runtime behavior. For example, detecting the semantic equality between the recursive and iterative implementation of the factorial computation. Semantic clone detection is the de facto technical boundary of clone detectors. This boundary was tested over the last years with interesting new approaches. This work contributes a semantic clone detection approach that detects clones with 0% syntactic similarity.

We present Semantic Clone Detection via Probabilistic Software Modeling (SCD-PSM) as a stable and precise solution to semantic clone detection. PSM builds a probabilistic model of a program that is capable of evaluating and generating runtime data. SCD-PSM leverages this model and its model elements to finding behaviorally equal model elements. This behavioral equality is then generalized to semantic equality of the original program elements. It uses the likelihood between model elements as a distance metric. Then, it employs the likelihood ratio significance test to decide whether this distance is significant, given a pre-specified and controllable false-positive rate. The output of SCD-PSM are pairs of program elements (i.e., methods), their distance, and a decision whether they are clones or not.

SCD-PSM yields excellent results with a Matthews Correlation Coefficient greater 0.9. These results are obtained on classical semantic clone detection problems such as detecting recursive and iterative versions of an algorithm, but also on complex problems used in coding competitions.

Index Terms—clone detection, semantic clone detection, probabilistic modeling, multivariate testing, software modeling, static code analysis, dynamic code analysis, runtime monitoring, inference, simulation, deep learning

I. INTRODUCTION

Copying and pasting source code fragments lead to code clones that are considered an anti-pattern. Code clones increase maintenance costs [1], [2], promote bad software design [3], [4], [5], and introduce or propagate bugs [6], [7], [8]. However, duplicating code fragments allows faster adaption to requirements, the reuse of stable and well-tested solutions [9], [10], and helps to overcome language limitations [11], [12] lowering development costs. The impact of code clones and the contradicting evidence various studies provide are the topics of an ongoing discussion in the community. However, developers will continue duplicating source code to leverage their benefits despite their drawbacks. Key is the awareness and management of clones to maximize efficiency while balancing quality.

Traditionally, the clone taxonomy distinguishes between four types of clones [12], [13], [14]. Type 1-3 describe code clones caused by copying and pasting the source code with or without changes. Type 4 clones describe code clones that do not have any syntactic similarity but implement the same functionality (semantic equivalence). For example, the recursive and iterative implementation of an algorithm (e.g., Fibonacci computation) have no syntactic similarity while implementing the same functionality. Existing tools have limited or no capabilities to detect Type 4 clones [15]. Most current studies exclude Type 4 clones because of the lack of tool support [16], [12], [13], [17], [18]. Nevertheless, Type 4 clones exist, and recent research efforts try to deepen the understanding of them [15], [19], [20]. This work provides a significant contribution to semantic clone detection in the form of novel concepts and a prototype implementing them.

We present *Semantic Clone Detection via Probabilistic Software Modeling (SCD-PSM)*. SCD-PSM extends on the Probabilistic Software Modeling (PSM) framework from our previous work [21] via a semantic clone detection pipeline. PSM builds probabilistic models (PMs) from programs. It analyzes the static structure and runtime behavior and replicates the program in the form of a generative probabilistic model. These models allow developers to reason about a program’s semantics via causal reasoning. SCD-PSM extends this work by leveraging the PMs and causal reasoning to find semantically (behaviorally) equivalent code elements (e.g., methods or properties with similar runtime behavior). SCD-PSM allows full quantification of the behavioral distance of code elements via likelihoods (i.e., probabilities). Furthermore, the likelihood evaluation via PMs allows for statistical significance tests (i.e., likelihood ratio test) to decide whether a pair of code elements are clones or not. SCD-PSM detects semantic clones with no textual similarity, such as the iterative and recursive version of an algorithm mentioned above. The most crucial feature of SCD-PSM is the active control of the false-positive rate, which often is an exclusion for the industry. This work extends our preliminary work on SCD-PSM where we introduced the notion of the model similarity [22].

Section II provides the background needed to understand SCD-PSM including the basics of PSM. Section III clarifies

what semantic clones are in the context of this work. Section IV presents the approach in which representation, the search space, and the various similarity stages are described. The study in Section V evaluates the concepts of the approach and Section VI discusses its results. Limitation of the approach and possible threats of validity of the study are given in Section VII and Section VIII. Section IX sets the results in context to the current state-of-art and Section X concludes this work.

```

1 int fa(int n){
2   product = 1
3   for(i = 1; i <= n i++)
4     product *= i
5   return product
6 }

```

Listing 1: A *for* implementation of factorial

```

1 int fb(int n){
2   product = 1
3   i = 1
4   while(i <= n)
5     product *= i
6     i++
7   return product
8 }

```

Listing 2: A *while* implementation of factorial

```

1 int fc(int n){
2   if(n <= 1) return 1
3   return fc(n - 1) * n
4 }

```

Listing 3: A *recursive* implementation of factorial

```

1 int fd(int n, String guard){
2   if(n < 1 && guard == "val") return -1
3   if(n < 1 && guard == "throw") throw Exception()
4   return fc(n)
5 }

```

Listing 4: A *delegate* implementation of factorial

II. BACKGROUND

The clone detection research community has a long history and defines many concepts, algorithms, and tools. In contrast, Probabilistic Software Modeling (PSM) is relatively new and combines software engineering and probabilistic modeling. Some terms need clarification; others require an introduction if they diverge from their traditional domain names.

A. Clone Detection

Clone detection is the process of finding two similar program fragments. Listings 1 to 4 are four different implementations of the factorial function ($n!$). Listing 1 is a *for*-loop implementation, Listing 2 uses a *while*-loop, and Listing 3 is recursively defined. Finally, Listing 4 delegates its implementation to $fc()$ from Listing 3 but may also return -1 in case of invalid inputs (including $n = 0$).

Clone detection revolves around the representation, pairing, and the similarity evaluation. *Representations* can be, e.g., text, graphs (e.g., AST), or probabilistic models like in this work. *Pairing* describes the selection of two fragments that

are potentially clones (e.g., $fa()$ and $fb()$). Each pair is called a *candidate clone pair* (or candidate pair). The *similarity evaluation* measures the similarity between the fragments of a candidate pair (e.g., counting the number of different characters). Finally, the *clone decision* labels the candidate pair as a clone or not given a criterion on the similarity (e.g., less than ten different characters).

The properties of the similarity metric split clones into two groups [12]. Type 1-3 clones capture textual similarity while Type 4 clones capture semantic similarity [13], [16], [23], [12], [14], [24].

Type 1: (Exact Clones) Program fragments that are identical except for variations in white-space and comments.

Type 2: (Parameterized Clones) Program fragments that are structurally/syntactically similar except for changes in identifiers, literals, types, and comments.

Type 3: (Near-Miss Clones) Program fragments that include insertions or deletions in addition to changes in identifiers, literals, types, and layouts.

Type 4: (Semantic Clones) Program fragments that are functionally/semantically similar (i.e., perform the same computation) without textual similarity.

These types are increasingly challenging to detect, with Type 4 being the most complex one. Note, that the definition of *Semantic Clones* is often relaxed where up-to 50% syntactic similarity of the code fragments is allowed (e.g., BigCloneBench [17]). However, we consider these clones as complex Type 3 clones (addition, deletions, reordering) and *not* as semantic clones. This means that semantic clones in the context of this work are clones with no syntactic similarity except for per-chance similarities (e.g., similar parameter names).

We will use $a \simeq b$ to denote that a is a clone of b . Furthermore, $a \not\simeq b$ denotes that a is not a clone of b .

B. Programs & Code Elements

PSM generalizes object-oriented terms to describe *code elements* in a program. Code elements are *types* T , *properties* Pr , and *executables* Ex that refer to, e.g., classes, fields, and methods in Java [25]. Additional code elements are *parameters* Pr and *results* Re of executables that refer to parameters (arguments) and return values of a method. Properties, parameters, and results are *atomic* code elements that have identifiable states at runtime. Types and executables are *compositional* elements that act as a collection of atomic elements.

Types *declare* properties and executables, capturing structural relationships. Executables have behavioral relationships that are categorized into *Inputs* (I) and *Outputs* (O). *Inputs* are *received parameters* Pa^I , *read properties* Pr^I , and *requested invocation results* Re^I . *Outputs* are *returned executable results* Re^O , *written properties* Pr^O , and *provided parameters* Pa^O .

We will denote atomic elements in lowercase, and compositional elements in bold-face lowercase, e.g., n and fa in Listing 1. Executable results are named after their executables, e.g., fa in Listing 1. $fc = \{n^{Pa,I}, fc^{Re,I}, fc^{Re,O}\}$ denotes the code elements of Listing 3. For readability sake, we will

omit the superscript classifiers if it is unambiguously possible e.g., $\mathbf{fa} = \{n, fa\}$. The subset of *inputs* is denoted by $\mathbf{fc}^{\mathcal{I}} = \{n^{Pa, \mathcal{I}}, fc^{Re, \mathcal{I}}\}$ and *outputs* by $\mathbf{fc}^{\mathcal{O}} = \{fc^{Re, \mathcal{O}}\}$. Finally, the set of all input and output combinations is given by

$$ex^{\mathcal{IO}} = \{(i, o) \in ex^{\mathcal{I}} \times ex^{\mathcal{O}}\}. \quad (1)$$

For example, $\mathbf{fd}^{\mathcal{IO}} = \{(n, fd), (guard, fd)\}$ describes the IO pairs of `fd()`.

C. Probabilistic Software Modeling

Probabilistic Software Modeling (PSM) [21] is a data-driven modeling paradigm that transforms a program into a probabilistic model. PSM extracts the structure and behavior of a program. The *structure* are the code elements of the program as described in Section II-B. This includes the actual elements, but also their call and structural hierarchy. The *behavior* is the runtime of the program, i.e., the actual data that is transformed by the program. The resulting PM and its model elements (i.e., random variables) is a copy of the original program with its code elements.

Model elements are the code elements in the context of a PM, and there is a one-to-one relationship between them. We will reuse the notation of code elements for model elements and augment it via probabilities. $P(x)$ denotes the probability distribution of variable x , e.g., $P_{\mathbf{fa}}(n)$ denotes the distribution of input parameter of the `fa`-method. $p(x)$ denotes the probability of a specific event of a variable, e.g., $p_{\mathbf{fa}}(n = 2)$. While notationally indistinguishable, the notation for model elements reasons about the behavior of code elements instead of their structural properties.

Model elements can generate observations (i.e., runtime events), but also evaluate their likelihood. *Generation* (or Sampling) draws, either marginally or conditionally, observations from a model element simulating the execution of the corresponding code element. For example, drawing 100 observations from $\mathbf{fa} \sim P_{\mathbf{fa}}(n, fa)$, i.e., values for $n^{\mathcal{I}}$ and $fa^{\mathcal{O}}$, simulates 100 program executions of this method. An example for *conditional generation* would be $\mathbf{fa}_{|n < 10} \sim P_{\mathbf{fa}}(fa | n < 10)$ that only draws observations where $n < 10$. *Evaluation* takes observations and evaluates their likelihood under a model element. For example, $P_{\mathbf{fa}}(n = 4, fa = 24)$ evaluates the likelihood of input 4 and output 24 under the `fa` model element. This extends to the underlying code element asking how likely is the method `fa` called with 4 while returning 24. Generation and evaluation are the core of any PSM applications and of SCD-PSM. How generation and evaluation work in detail and from a technical perspective is given by Thaller et al. [26].

III. SEMANTIC CLONES

A clear understanding of what SCD-PSM defines as *semantic clone* is essential in understanding the approach and its design choices.

Definition 1. *A semantic clone is a pair of executables whose (partial) input, and output relationships exhibit significant (conditional) similarities.*

Definition 1 defines semantic clones over the similarity between IO relationships of executables. This holds if the IO relationships are only partially similar, i.e., not all combinations of IO pairs between executables have to be similar. For example, `fd` in Listing 4 has two IO pairs ($\mathbf{fd}^{\mathcal{IO}} = \{(n, fd), (guard, fd)\}$) while `fa` in Listing 1 has one IO pair ($\mathbf{fa}^{\mathcal{IO}} = \{(n, fa)\}$). According to the definition, at least one IO pair comparison needs to be similar such that both executables are declared as a semantic clone (e.g., $(n, fd) \simeq (n, fa)$).

Furthermore, the similarities between IO pairs may only be conditional, i.e., the similarity of matching IO pairs might be suspect to the state of any other code element in the comparison context. For example, the IO pair $(n, fd) \simeq (n, fa)$ is only a perfect clone in case that `fd.guard != "val"`. If `fd.guard == "val"` the IO behavior would differ in case of $n = 1$ (`fd(1) ↦ -1` while `fa(1) ↦ 1`). According to the definition, at least parts of the behavior need to be similar, capturing complex multidimensional behavioral patterns in IO relationships.

The rationale behind the comparison of IO relationships is one of cause and effect. If a pair of executables exhibit similar effects given similar causes, then their computational behavior is identical. Extending this rationale by multiple inputs and outputs results in multiple IO relationships leading to *partially conditional similarity*.

IV. APPROACH

Figure 1 illustrates SCD-PSM. It is a five-fold approach in which:

- A. **[Modeling]** PSM builds a probabilistic model that reflects the original program;
- B. **[Search Space]** A search space of candidate pairs is constructed by pairing executable model elements;
- C. **[Static Similarity]** The static similarity accepts candidate pairs with matching data types;
- D. **[Dynamic Similarity]** The dynamic similarity accepts candidate pairs with similar runtime data;
- E. **[Model Similarity]** The model similarity accepts candidate pairs with similar model behavior.

The approach represents a rejecting filter pipeline that candidate pairs must traverse to be declared a clone. Static, Dynamic, and Model Similarity represent filter stages of increasing complexity.

The main contribution of this work is the implementation of a semantic clone detection pipeline on top of PSM. Further, we provide a simple yet effective process of traversing the search space of potential clones while respecting the pipeline's performance properties. Finally, we show that the behavioral equivalence of model elements generalizes to the semantic equivalence of code elements.

A. Modeling

Starting from the *Source Code* in Figure 1, PSM builds a *Probabilistic Model* (PM) [21] of the program (1). SCD-PSM uses this PM as a representation for the clone detection similar to text-based clone detectors that use text fragments. While

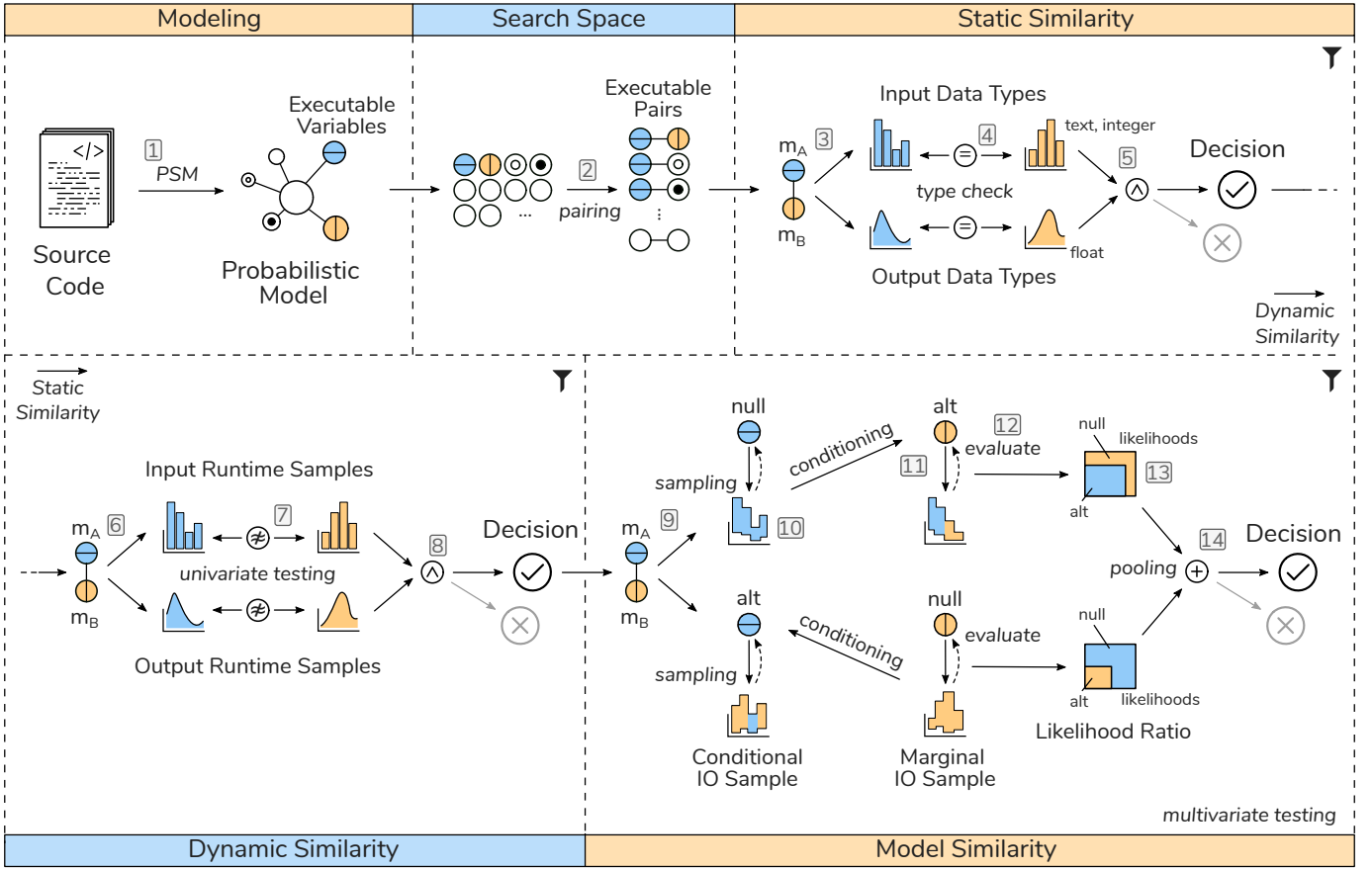


Figure 1: The modeling phase transforms the program into a PM. The search space phase then pairs the PM model elements into candidate pairs. Finally, Static-, Dynamic- and Model Similarity evaluates the behavioral equality of the candidates.

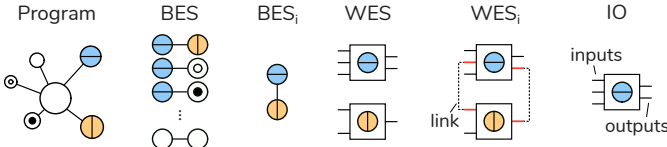


Figure 2: SCD-PSM operates on four levels of abstraction: program, between executable, within executable, and the IO level.

being the first step of SCD-PSM, we consider the PM as given since we only use PSM to build a model.

Executable model elements in the PM act as a surrogate to the executables in the original program (e.g., $f_a()$). SCD-PSM uses these elements and computes their similarity finding behaviorally equivalent model element pairs and extrapolates this knowledge to semantically equivalent code element pairs. In conclusion, the SCD-PSM allows for method-level semantic clone detection based on PMs representing the original executables in the program.

B. Search Space

SCD-PSM conducts method-level semantic clone detection, which operates on multiple abstraction levels. Figure 2 illus-

trates these levels starting with the program and ending with the inputs and outputs of an executable.

The second step in Figure 1 builds a *within- and between-executable space* that SCD-PSM searches for clones. The *Between-Executable Space (BES)* is the set of executable combinations

$$BES = \{\{a, b\} \in Ex \times Ex \mid a \neq b\}, \quad (2)$$

where exa, exb is a *candidate pair* (or executable pair), and Ex is the set of all executables in the current analysis (illustrated in Figure 2). The theoretical size of the between-executable space are all 2-length combinations without replacement given by

$$|BES| = \frac{|Ex|!}{2 \cdot (|Ex| - 2)!}, \quad (3)$$

where $|\cdot|$ describes the size of the underlying set. Note that the size of the BES is smaller than the cartesian product since $\{a, b\} = \{b, a\}$. Figure 1 shows this pairing process in the Search Space aspect (2). The *Within-Executable Space (WES)* is the product of IO pairs

$$WES^{ab} = \{(i, j) \in a^{IO} \times b^{IO}\}. \quad (4)$$

Figure 2 illustrates the WES and one IO pair from the WES that

we also call *link*. The theoretical size of the within-executable space is

$$|\mathbf{WES}^{ab}| = |\mathbf{a}^{\mathcal{IO}}| \cdot |\mathbf{b}^{\mathcal{IO}}| \quad (5)$$

For the sake of visualization IO pairs are not shown in Figure 1 but are abstracted in their executable elements. The maximum theoretical search space is

$$S = \sum_i |\text{wes}(\mathbf{BES}_i)|, \quad (6)$$

given that *wes* describes a construction function according to Equation (4), and \mathbf{BES}_i is the i 'th candidate pair.

In practice, SCD-PSM evaluates only a fraction of possible combinations because of the skip evaluation. The *skip evaluation* consists of two search space limiting factors: greedy evaluation and transitive similarity. *Greedy evaluation* stops the search through the WES once a similar pair is found. The initial detection process is only interested in whether there is a similarity or not. Post-analysis can then find all possible similarities for the program comprehension and developer guidance. *Transitive similarity* skips evaluations in the BES because if $a \simeq b \simeq c$ then also $a \simeq c$ holds. In sum, *skip evaluation* allows for efficient processing by stopping the search through the WES once a clone pair is found, and by skipping pairs in the BES if the transitive similarity is given.

In conclusion, SCD-PSM compares IO pairs of executable model elements and uses greedy evaluation and transitive similarity to traverse the search space efficiently.

C. Static Similarity

The static similarity stage is a filter that accepts candidate pairs based on their data type, as shown in Figure 1. Data types in a PSM model are integers, floats, and text.

Input (3) of the stage are the IO pairs $\mathbf{WES}^{ab} = \text{wes}(\{\mathbf{a}, \mathbf{b}\})$ of a candidate. The filter *criteria* (4) accepts a candidate pair if at least *one* link (i.e., IO pair) has a matching data type, i.e., the input but also the output have a matching data type. *Output* (5) is a boolean decision whether the candidate pair is a clone or not from a static viewpoint. If positive, then the candidate pair is moved to the next pipeline stage, i.e., the *Dynamic Similarity* evaluation (see Figure 1). If negative, then the candidate pair is marked as being *not* a clone $\mathbf{a} \not\simeq \mathbf{b}$ and no further processing is conducted. For example, the IO pairs $(n, fa) \simeq (n, fb)$ would be statically accepted as clones as both inputs and outputs have the same data type (integer). A counterexample is given by $(n, fa) \simeq (guard, fd)$ where the input data types are integers and text.

The current version of the prototype does not employ any complex static analysis to detect, e.g., unaltered delegations as given in `fd`. However, this optimization can be added to the pipeline. In conclusion, the static similarity stage filters candidates based on their data type.

D. Dynamic Similarity

The dynamic similarity stage is a filter that accepts candidate pairs based on the runtime data, as shown in Figure 1.

Input (6) of the stage are again the IO pairs of a candidate as in the static similarity Section IV-C. The filter *criteria*

(7) accepts a candidate pair if at least *one* IO pair has an *insignificant* diverging runtime distribution. The divergence of inputs and outputs can be computed via univariate statistical tests. This work uses the Kolmogorov-Smirnov test [27], which is a general non-parametric test. *Output* (8) is a boolean decision whether the candidate pair is a clone or not from a dynamic viewpoint. Again, the candidate is removed from the pipeline or forwarded to the next pipeline stage, depending on the decision. For example, the IO pair $(n, fa) \simeq (n, fd)$ with `guard == true` fails the filter given that runtime events with $n = 0$ reach a critical mass. In comparison, $(n, fa) \simeq (n, fb)$ would be accepted by the stage.

A requirement is that the program uses a synthetic trigger. Otherwise, the comparison of the data distributions may fail because of the different *modus operandi* of the program. For example, running `fa` and `fb` where $n_{fa} = \mathcal{U}(0, 4)$ and $n_{fb} = \mathcal{U}(5, 10)$ would cause the dynamic stage to fail even if the implementations are equivalent. Property-based testing [28] can be used for this purpose to generate diverse synthetic inputs.

In conclusion, the dynamic similarity pre-filters candidates based on univariate tests on the input and output events.

E. Model Similarity

The model similarity stage is a filter that accepts candidate pairs based on the model data, as shown in Figure 1. The model data is sampled from the models of the executables in the candidate pairs. This stage conducts a multivariate test as opposed to the univariate test the dynamic similarity computes. The multivariate testing considers the conditional similarities that are not part of the currently inspected IO behavior but maybe conditionally influencing the IO pair. For example, $(n, fd) \simeq (n, fa)$ holds but is conditionally dependent on *guard*. The dynamic stage can only look at the marginal (i.e., average) behavior of the IO pair, while the model similarity can factor *guard* into its decision.

Input (9) are the IO pairs of a candidate $\mathbf{WES}^{ab} = \text{wes}(\{\mathbf{a}, \mathbf{b}\})$. Then the cross-wise log-likelihood ratio of the models is computed by (*conditional*) *generation* and *evaluation*. *Output* is a boolean decision whether the candidate pair is a clone or not from a model viewpoint. Figure 1 illustrates the entire process of the model similarity.

- 1) A reference model $M^{null} = \mathbf{a}$ and an alternative model $M^{alt} = \mathbf{b}$ (candidate pair executables).
- 2) A IO-pair $p = \mathbf{WES}_i^{ab}$ is selected that is target of the comparison (link).
- 3) A reference sample D^{null} is drawn from M^{null} (10).
- 4) A alternative sample $D^{alt|null}$ is drawn from M^{alt} by optimizing towards the p dimensions in the D^{null} effectively conditioning the drawn samples (11).
- 5) D^{null} is evaluated under M^{null} resulting the reference log-likelihood LL^{null}
- 6) $D^{alt|null}$ is evaluated under M^{alt} (12) yielding the alternative log-likelihood LL^{null} .
- 7) Finally, likelihood ratio of the link is computed

$$\lambda = LL^{alt} - LL^{null} \quad (7)$$

This procedure is repeated for the reverse assignment of null and alt roles. Finally, both log-likelihood ratios are combined by a pooling operator into the final decision (14).

The role-swap is needed to avoid sub-model relationships. For example, if $M^{null} = \mathcal{N}(0, 3)$ and $M^{alt} = \mathcal{N}(0, 1)$ then LL^{alt} will be very high because M^{alt} is a sub-model from M^{null} . However, reversing the roles highlights the differences in the models.

The final decision is based on the Generalized Likelihood Ratio Test (GLRT) [29]. It measures whether the log-likelihoods are significantly different from 0 where λ (see Equation (7)) is the test statistic. The null hypothesis is that the models are significantly different and rejects for small ratios $\lambda \leq c$ where c is set to an appropriate Type 1 error (false-positive rate). For example, $\lambda < \log(0.01)$ allows 1 out of 100 candidates to be a false-positive, i.e., wrongly rejecting semantic equivalence. The pooling operator combines the link results either via hard or soft pooling. *Hard pooling* conducts for both links a GLRT yielding a positive decision if *both* links are positive. *Soft pooling* averages the link log-likelihoods ratios and then computes the GLRT yielding a positive decision if the joint GLRT is positive. Hard pooling does not allow any sub-model relationships, while soft pooling relaxed this constraint.

In conclusion, the model similarity conducts a multivariate significance test between two models, including possible conditional dependencies.

V. STUDY

This study answers the following research questions.

- RQ1 Does behavioral equality between model elements generalize to semantic equality of code elements?
- RQ2 Does the skip evaluation significantly improve on the computation demand of SCD-PSM?
- RQ3 Does the skip evaluation negatively impact the detection performance (i.e., precision, recall, and MCC)?

RQ1 answers the essential question if semantic clones can be detected via SCD-PSM. RQ2 answers the question of whether the search space can be efficiently processed using skip evaluation. RQ3 answers the question of how the skip evaluation influences the performance of the detection process. This is important because candidate pairs might be skipped based on false-positives or false-negatives.

A. Setup

We implemented a prototype for SCD-PSM on top of Gradient [21], a prototype for PSM. The elements and data flow of the detection process are given by Figures 1 and 2.

- 1) The input *Source Code* were 13 different clone classes with a total of 108 implementation variants. This includes classical algorithms implemented recursively and iteratively (e.g., bubble sort), but also hard problems from the programming competition Google Code Jam.
- 2) The *Probabilistic Model* was computed via Gradient, a PSM prototype. We used the same hyper-parameters as reported in our previous work [21].
- 3) The *Search Space*, i.e., the *BES* and *WES*, was created according to Section IV-B.

Table I: Semantic clone detection dataset consisting of 108 instances across 13 clone classes.

Problem	Type	Property	Source Code			Clone	
			Executable	Parameter	Loc	Class	Instance
Factorial	2	0	4	2	21	A	2
Fibonacci	2	0	4	2	10	B	2
Sort	4	0	12	18	110	C	4
R0AA	11	4	39	54	384	D	10
R0AB	10	0	33	27	374	E	10
R0AC	12	4	32	45	433	F	10
R1AA	10	0	26	41	543	G	10
R1AB	13	8	28	55	595	H	10
R1BA	12	4	31	53	269	I	10
R1CA	17	18	35	62	463	J	10
R1CB	15	18	34	82	772	K	10
R1CC	10	0	21	41	338	L	10
R2AA	10	0	23	18	471	M	10
	128	56	322	522	4783	13	108

- 4) Each valid candidate pair was then flown through the *Static-*, *Dynamic*, *-* and *Model-Similarity* and filtered according to Sections IV-C to IV-E. Candidates that passed the entire filter pipeline were marked as clones.

We apply NiCad [30], PMD [31], and iClones [32] to the dataset to provide some context of our results and the quality of the dataset. NiCad is a text-based, PMD a token-based, and iClones a graph-based clone detector. Each of them can detect Type 1-3 clones. We were not able to run any of the existing clone detectors that report Type 4 capabilities (Tisem [33], AutoCodeGEN [34], DeepSIM [35]) because of technical or availability issues.

All experiments were computed on a single machine (Intel i7).

B. Subject Examples

The study uses three well-known algorithms and 10 Google Code Jam 2017 (GCJ) ¹ problems as listed in Table I. The total dataset contains 108 implementation variants across 13 clone classes described by *Instance*. Factorial, Fibonacci, and Sort do not need any further explanation. The GCJ problems are well specified complex optimization problems packaged in an everyday theme. R0AA in Table I is the round 0 (qualification), category A, and task A problem. It asks developers to compute the minimum amount of flips of pancakes that are needed for all pancakes to be on the upside. The constraint is that always k pancakes are flipped at the same time.

The dataset contains in total 5778 (see Equation (3)) candidate pairs of which 458 are semantic clones and 5320 are not. This yields a positive to negative ratio of 1/11.6, indicating a highly imbalanced distribution. A similar or even more pronounced imbalance is to be expected in real-world applications.

Each instance was triggered with input data to allow PSM to model the different implementations. Factorial, Fibonacci, and Sort were triggered by sampling from a uniform distribution

¹<https://codingcompetitions.withgoogle.com/codejam/archive>

$U(0, 20)$. GCJ problems were triggered by the input data provided by the competition. Each instance received the same trigger.

GCJ problems read and write to the standard stream, which is impractical in terms of reproducibility. Our dataset is constructed such that each implementation has a `run`-method representing the cloned executable. The study results are limited to the `run`-method even if the solutions use auxiliary methods. This guarantees a proper problem scope, well-defined recall and precision values, and a clearly defined benchmark for future reproducibility. The dataset is openly available ².

C. Controlled Variables

The study controls for the search space *Evaluation* strategy, *Dynamic False-Positive Rate (D-FPR)*, *Model False-Positive Rate (M-FPR)*, and *Pooling*.

Evaluation describes how the search space is processed: *exhaustive*, or *skip*. The exhaustive evaluation compares each executable candidate with each other. The skip evaluation uses the transitive similarity (see Section IV-B) and may skip evaluation if possible.

Dynamic False-Positive Rate (D-FPR) defines the critical value α of the Kolmogorov-Smirnov test with 0.001 and 0.01 at which similarity is rejected.

Model False-Positive Rate (M-FPR) defines the critical value c of the Generalized Likelihood Ratio test with 0.001 and 0.01 at which similarity is rejected.

Pooling defines how the likelihood ratios from the two link directions are combined (see Figure 1, (8)) with values: *hard*, or *soft*. *Hard* pooling evaluates whether each link reaches the critical value c and accepts the clone if both links evaluate positive.

$$\lambda_{Link_A} \leq \frac{\log c}{2} \text{ and } \lambda_{Link_B} \leq \frac{\log c}{2} \quad (8)$$

Soft pooling evaluates the average log-likelihood ratios (geometric mean of likelihoods)

$$\frac{\lambda_{Link_A} + \lambda_{Link_B}}{2} \leq \log c, \quad (9)$$

and compares it against the critical value c .

An additional fixed parameter is the *number of particles*. *Number of particles* defines the samples size that is generated during the model similarity $|D| = 50$.

D. Response Variables

The response measures of the study are the number of *Skip Evaluations*, processing *Duration*, *TP*, *FP*, *TN*, *FN*, *Precision*, *Recall*, *Accuracy*, and *Matthews Correlation Coefficient*.

Skip Evaluations measures the number of evaluations that were skipped because of the skip evaluation strategy.

Duration measures the elapsed time to compute one candidate pair.

TP, FP, TN, FN measures the True Positive (TP), False Positive (FP), True Negative (TN), and False Negative (FN) detection results compared to the ground truth.

Precision measures the fraction of detected clones that are truly clones

$$\frac{TP}{TP + FP} \quad (10)$$

Recall measures the fraction of semantic clones that have been found

$$\frac{TP}{TP + FN} \quad (11)$$

Accuracy measures the closeness of the detector to the ground truth

$$\frac{TP + TN}{TP + FP + TN + FN} \quad (12)$$

Matthews Correlation Coefficient (MCC) measures the quality of the clone detection in form of a correlation ranging from -1 to 1 with 0 being a random selection. MCC is robust against class imbalance because it considers the entire confusion matrix. The MCC will be the reference performance metric in evaluating the success of the approach using the Evans [36] interpretation guidelines.

$$\frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (13)$$

E. Experiment Results

Creating the PSM model with Gradient took 2134.38 s, resulting in an average modeling time of 19.75 s for the 108 executables.

Table II contains the aggregate results of the top-5 experiments along with the results of the worst experiment. The bottom line in Table II is the average performance of the top-5 experiments. The generally expected performance of the approach is an MCC of 0.965. High confidence for negative examples is given with no false-positives reflecting the pipelines FPR rates (D-FPR \times M-FPR).

A multiple linear regression was calculated to predict MCC based on the control variables: evaluation, D-FPR, M-FPR, and pooling. A significant regression equation was found ($F(4, 11) = 39.8, p < 0.01$), with an R of 0.911. MCC increased by 0.056 for the *skip* evaluation and by 0.080 and for *soft* pooling. Larger D-FPR and M-FPR decrease the MCC by -0.042 and -5.455 . All predictors were significant ($p < 0.01$) except for D-FPR.

The best experiment is given with a *skip evaluation*, *0.100* D-FPR and *0.001* M-FPR rates, and *soft pooling* (Nr. 1) with an MCC of 0.975. A total of 345 candidates were skipped while reaching a recall of 0.933. Table III lists the cumulative performance of the best model, starting with an initial prediction that all candidates are semantic clones (rejecting pipeline). The *static* stage finds 71.729 % (3816) of the FPs improving the MCC by 0.409. The *dynamic* stage additionally removes another 27.330 % (1454) of FPs but introduces 1.528 % (7) of the possible FNs. An improvement of the MCC by 0.527 is achieved via the dynamic stage. Finally, the *model* stage removes the remaining 0.939 % (50) FPs but introduces additional 3.056 % (14) FNs. The model stage improves the MCC by 0.039.

²github.com/gradient-dataset-clones

Table II: Results and the average result of the top-5 of 16 experiments along with worst experiments.

Nr	Controlled Variables				Response Variables										
	Evaluation	D-FPR	M-FPR	Pooling	Duration	TP	FP	TN	FN	Skip	Accuracy	Precision	Recall	F1	MCC
1	skip	0.100	0.001	soft	28	437	0	5320	21	345	0.996	1.000	0.954	0.977	0.975
2	skip	0.010	0.001	soft	27	437	0	5320	21	345	0.996	1.000	0.954	0.977	0.975
3	exhaustive	0.010	0.001	soft	28	425	0	5320	33	0	0.994	1.000	0.928	0.963	0.960
4	skip	0.010	0.010	soft	32	423	0	5320	35	332	0.994	1.000	0.924	0.960	0.958
5	exhaustive	0.100	0.001	soft	34	421	0	5320	37	0	0.994	1.000	0.919	0.958	0.955
16	exhaustive	0.100	0.010	hard	47	293	0	5320	165	0	0.971	1.000	0.639	0.780	0.787
1-5	skip	0.010	0.001	soft	29	428	0	5320	29	340	0.995	1.000	0.936	0.967	0.965

Duration in minutes

Table III: Performance breakdown of the best performing experiment listed as Nr. 1 in Table II.

Nr	Stage	TP	FP	TN	FN	Accuracy	Precision	Recall	F1	MCC
1	initial	458	5320	0	0	0.079	0.079	1.000	0.147	–
2	static	458	1504	3816	0	0.740	0.233	1.000	0.379	0.409
3	dynamic	451	50	5270	7	0.990	0.900	0.985	0.941	0.936
4	model	437	0	5320	21	0.996	1.000	0.954	0.977	0.975
		437	0	5320	21	0.996	1.000	0.954	0.977	0.975

The worst experiment is given with *exhaustive evaluation*, D-FPR of 0.100, M-FPR of 0.010, and *hard pooling* (Nr.16) with an MCC of 0.787.

On average 5.884% (340) of the total 5778 evaluations could be skipped. This equals 74.235% of the total 458 TPs. On average 0.626% (50354) of the total 134782 IO pair evaluations could be saved via greedy evaluation. The average duration of the exhaustive experiments was 2394s leading to 414ms per candidate. Skip experiments lasted on average for 1988s with 344ms per candidate. The static stage lasted on average for 100.744 μ s (<0.001%), the dynamic stage for 208.386ms (0.106%), and the model stage for 1.749s (0.893%) per candidate. Each candidate may contain multiple IO pairs (WES) that need to be evaluated.

Table IV: Results of PMD, NiCAD, and iClones on the semantic clone detection dataset.

Tool	TP	FP	TN	FN	Precision	Recall	F1	MCC
PMD	0	0	5320	458	–	0.000	0.000	–
NiCAD	0	0	5320	458	–	0.000	0.000	–
iClones	1	0	5320	457	1.000	0.002	0.004	0.045

Table IV lists the detection results of PMD, NiCAD, and iClones. In total, one clone was detected by iClones containing a complex matrix computation. All the other tools failed to find clones in the dataset.

VI. DISCUSSION

The goal of the study was to provide evidence of whether behavioral equality of model elements generalizes to semantic equality of code elements (RQ1). Furthermore, we were interested in the skip evaluation and its performance implications (RQ2 and RQ3).

Table II and Table III present strong results in favor for RQ1. The MCC for the top-5 experiments was *very strong* with all MCCs being above 0.9. Even the worst experiment still yielded a *moderate* performance of 0.749. Table IV brings these numbers into a context where only *one* tool could detect a clone. These detectors are well-known for their capability to detect Type 1-3 clones indicating the quality of the dataset. The solutions in the dataset itself contain differential tests confirming the semantic equivalence. In conclusion, given the strong evidence provided by the performance of SCD-PSM and the lack of performance of our comparison tools, we conclude that RQ1 holds.

RQ1 – Behavioral equality between model elements generalizes to semantic equality of code elements allowing for semantic clone detection via probabilistic software modeling.

The goal of the static and dynamic stage is to reduce the number of evaluations that the model stage must conduct. Each stage incurs an increased evaluation cost per candidate, with the model stage taking the largest share of 89% of the evaluation time. Every TP has to pass the model stage to be declared a clone (rejecting pipeline). The skip evaluation avoided, on average, the recomputation of 74% (340) of the TP candidate pairs. This offloads most of the evaluation time to the earlier stages that are computationally inexpensive. In sum, the skip evaluation reduces the number of model evaluations that are responsible for most of the evaluation time to a quarter.

RQ2 – Skip evaluation reduces the number of evaluations for the most expensive stage (model) in the SCD-PSM pipeline significantly.

The current skip evaluation has further optimization potential. In the study, we randomized the candidate pairs fed to the pipeline. The randomization avoids a performance bias caused by the dataset structure where clone classes are organized in packages (e.g., `r0AA.Dev0.run(...)`). While opportunities to skip an evaluation are randomly encountered in our prototype, greedily exploring a clone class until all possible candidates are exhausted would be a better search strategy. Another possibility is to pre-compute the dynamic similarity stage for all candidates, which provides a distance metric between executables (i.e., the discrepancy of the KS-test). This distance can then be used to sort candidates such that clone classes are clustered and ranked earlier in the evaluation process leading to more skip evaluations. The last and most impactful performance optimization is incremental clone detection, which is naturally supported by SCD-PSM. Each model represents a snapshot of an executable in a specific version. Hence, comparing across versions means that only new and modified executables need to be examined.

Skip evaluation can cause cascading errors given an FP. Once an FP is introduced, every semantic clone to the FP has a chance to become an FP in the same clone classes itself. These cascading FPs are potential sources of significant performance degradations. Skip evaluation experiments are ranked higher and are significantly better than experiments that conducted an exhaustive search. However, the absolute performance gain is merely an MCC of 0.056 hinting a performance significance introduced by the small sample size (16 experiments). Nevertheless, given the evidence in Table II, and Section V-E we can conclude that skip evaluation does not affect the performance of the detector.

RQ3 – The skip evaluation has no negative impact on the performance of the detector given low false-positive rates.

One last aspect worth a discussion are the application pragmatics of SCD-PSM. Creating the PSM model for the 108 executables takes about 40 min. The clone detection for these executables lasts 29 min. This means the SCD-PSM is not intended to be used within an Integrated Development Environment (IDE) as live detection mechanism, but rather in a Continuous Integration (CI) pipeline triggered at specific versions. To bring these durations into a better context, the projects Structurizr, JLatexmath, and PMD had 50 to 574 models [21]. Nevertheless, limiting the number of analyzed executables to a sensible subset of the system is still essential.

In conclusion, SCD-PSM can detect semantic clones with high precision and recall.

VII. LIMITATIONS

SCD-PSM inherits the limitations of PSM.

PSM does only model data in an application which exclude reference data types. References are every changing addresses with no relation to the running program. Hence, they have no meaningful underlying distribution that can be modeled. However, once references are dereferenced, e.g., by accessing a field, their accessed data will be part of the model. Once part of the model, it can be used in the context of SCD-PSM. Nevertheless, algorithms with the sole purpose of manipulating and comparing references are unfit for SCD-PSM.

PSM explodes lists into singular values since distributions do not contain any order information. This means, executables that change the order of sequences are matched based on the values, not their order. As a consequence, invoking a wrong implementation, e.g., of a sorting algorithm, may lead to false positives. Extending PSM to distributions of sequences alleviates the issue but is not a trivial task.

Another limiting factor is that the approach needs to run the program to build the model. This means that SCD-PSM can only be applied to runnable programs.

```
1 int fx(int n){
2   product = 1
3   for(i = 1; i <= n i++)
4     product /= i
5   return product
6 }
```

Listing 5: Changed *for* implementation of factorial

Finally, SCD-PSM cannot detect Type 2-3 clones since textual similarities represent a different problem set. A proof can be easily constructed by applying a textual detector and a semantic clone detector on Listings 1 and 5. Listing 5 changes the multiplication symbol to a division symbol, which most textual detector will ignore reporting $fa \simeq fx$. However, the behavior of the method completely changed, and a semantic clone detector would report $fa \not\approx fx$. Inversely, we can add a random number of syntactically correct statements to a perfect clone of Listing 1. These statements may be added in any line of the methods (e.g., $1 * 2 + 2 - 1$) as long as it does not change the behavior. Textual clone detectors will at some point report $fa \not\approx fx$ because it contains too many additional statements. However, semantic clone detectors will report $fa \simeq fx$ as long as the added statements are neutral to the original behavior.

VIII. THREATS TO VALIDITY

An external threat to validity is given by the type of problems in the dataset. The dataset problems are complex optimization problems that are not encountered in this density in real-world applications. Applying SCD-PSM to, e.g., a web-application might incur a change of performance. While a valid criticism, the alternative would be to detect clones in random real-world applications opportunistically. However, this approach lacks a proper notion of recall and precision and needs additional human inspections, which are prone to inconsistencies [37], [24]. We valued correct performance

metrics on a stable and clearly defined benchmark higher than the potential performance loss or gain that arises from a problem context switch. Furthermore, real-world applications will tend to be less algorithmic intensive, thus potentially easier to detect.

IX. RELATED WORK

We started this work by defining what *semantic clones* means in the context of our approach (Section III). While our definition is motivated in the capabilities of our approach, we can see strong similarities to the definition of Juergens [15]. Both definitions define behavioral similarity via IO relationships. Also, Juergens already discussed a notion of partial and conditional similarity. This understanding of Type 4 clones can be seen in multiple more recent studies [38], [39]. In that, we see the progress of the community in terms of Type 4 clones as the definition becomes more specific.

Many studies evaluated textual clones. However, only a few studies are reporting results on semantic clones without relaxing the definition of Type 4. Rattan [14] et al. provided a review of clone detection studies. The review also investigated approaches that tackle Type 4 clones. They conclude that some approaches solve approximations (i.e., complex Type 3 clones) of Type 4 clones.

Horwitz [40] detected textual and semantic differences in programs via a Program Representation Graph, which is similar to a Program Dependency Graph (PDG). A PDG graph contains nodes for statements and control points. Edges represent data and control dependencies between these. PDG-based approaches [23], [41], [42] use these data and control dependencies to find similar sub-graphs between the candidates. PDG based methods can detect complex Type 3 clones, e.g., Listings 1 and 2. However, the compared PDG sub-graphs are a representation of the source code; thereby, the approaches still rely on syntactic similarity [43].

Test-based methods randomly trigger the execution of candidates and measure whether equal inputs cause similar outputs. Jiang and Su [44] were able to find semantically equivalent methods without any syntactical similarities. A similar approach was presented by Deissenboeck et al. [38]. One issue with test-based clone detection is that candidates need a similar signature. Differences in data types or the number of parameters can not be effectively handled by the test-case generators or the similarity measurement. SCD-PSM works similarly to test-based methods in that it observes the runtime and compares the resulting behavior. However, SCD-PSM builds generative models from the observed behavior capable of generating and evaluating data. Missing dimensions are imputed by conditioning and sampling via the models. This conditioned sample is then evaluated, allowing SCD-PSM to overcome the issue of signature mismatches. Furthermore, PSM abstracts the data types into text, integer, and floats mitigating data type mismatches.

Zhao and Huang [35] developed DeepSim, which learn control and data flow encodings that are then used in a binary classification task. DeepSim uses neural networks to learn a latent representation of the control and data flow. PSM also uses

neural networks (each model element) but learns an underlying representation of the runtime. Hence, the most significant difference is in *what* is modeled (static vs. runtime information). DeepSim also evaluated on a Google Code Jam dataset however used the problems from 2016. DeepSim reached an F1 score of 0.76 on the GCJ 2016 competition, while SCD-PSM reached 0.967 on the GCJ 2017. While not entirely comparable, we expect these results to be a good approximation since the type of problems (complicated optimization), and its context (competition with one file solutions) match.

Finally, the community still lacks in terms of datasets for semantic clones. We want to caution opportunistic clone detection datasets that randomly select real-world applications and try to find clones; especially, if the dataset is used to evaluate a new detection approach. Neither, the number of positive nor negative clone pairs is known. Often quantities are estimated via manually inspected or using multiple tools as reference. However, seldom the standard error of these measurements are provided. The high imbalance between clones and non-clones in real-world applications increases the risk of misestimating even further as a multiplicative factor.

X. CONCLUSIONS AND FUTURE WORK

In this work, we presented Semantic Clone Detection via Probabilistic Software Modeling (SCD-PSM). PSM builds a Probabilistic Model (PM) from a program that can be used to generate or evaluate runtime events. We used these PMs to detect semantic clones in programs that have 0% syntactic similarities.

We have discussed the representation, search space, static-, dynamic-, and model-similarity that form the main aspects of SCD-PSM. Our study evaluated SCD-PSM in great detail, with an average MCC greater 0.9. Also, the study showed the capability to control the false-positive rate, which is a significant point for industry adoption. Finally, we concluded that behavioral equality of model elements generalizes to semantic equality of code elements.

Our future endeavors will focus on defining and declaring a comprehensive semantic clone detection benchmark. The focus of this benchmark should lie on the sheer amount of clones, but rather the complex intricacies of the semantic behavior of programs. As a logical conclusion, we further want to extend the research on SCD-PSM by leveraging this benchmark.

In conclusion, SCD-PSM is capable of detecting semantic clones with 0% syntactic similarity. SCD-PSM clone decisions have a statistical foundation that allows to control for the false-positive rate. Finally, these decisions can be ranked based on the likelihood of the candidate being a semantic clone.

ACKNOWLEDGMENTS

The research reported in this paper has been supported by the Austrian Ministry for Transport, Innovation and Technology, the Federal Ministry of Science, Research and Economy, and the Province of Upper Austria in the frame of the COMET center SCCH.

REFERENCES

- [1] Mayrand, Leblanc, and Merlo, "Experiment on the automatic detection of function clones in a software system using metrics," in *Proceedings of International Conference on Software Maintenance ICSM-96*. Monterey, CA, USA: IEEE, 1996, pp. 244–253.
- [2] A. Monden, D. Nakae, T. Kamiya, S. Sato, and K. Matsumoto, "Software quality analysis by code clones in industrial legacy software," in *Proceedings Eighth IEEE Symposium on Software Metrics*. Ottawa, Ont., Canada: IEEE Comput. Soc, 2002, pp. 87–94.
- [3] R. C. Martin, Ed., *Clean Code: A Handbook of Agile Software Craftsmanship*. Upper Saddle River, NJ: Prentice Hall, 2009.
- [4] M. Fowler and K. Beck, *Refactoring: Improving the Design of Existing Code*, ser. The Addison-Wesley Object Technology Series. Reading, MA: Addison-Wesley, 1999.
- [5] A. Hunt and D. Thomas, *The Pragmatic Programmer: From Journeyman to Master*. Reading, Mass: Addison-Wesley, 2000.
- [6] A. Chou, J. Yang, B. Chelf, S. Hallem, and D. Engler, "An empirical study of operating systems errors," *ACM SIGOPS Operating Systems Review*, vol. 35, no. 5, p. 73, Dec. 2001.
- [7] Z. Li, S. Lu, S. Myagmar, and Y. Zhou, "CP-Miner: Finding Copy-Paste and Related Bugs in Large-Scale Software Code," *IEEE Transactions on Software Engineering*, vol. 32, no. 3, pp. 176–192, 2006.
- [8] R. Geiger, B. Fluri, H. C. Gall, and M. Pinzger, "Relation of Code Clones and Change Couplings," in *Fundamental Approaches to Software Engineering*, L. Baresi and R. Heckel, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, vol. 3922, pp. 411–425.
- [9] J. Krinke, "Is Cloned Code More Stable than Non-Cloned Code?" *Proceedings - 8th IEEE International Working Conference on Source Code Analysis and Manipulation, SCAM 2008*, pp. 57–66, 2008.
- [10] —, "Is Cloned Code Older than Non-Cloned Code?" 2011.
- [11] C. J. Kapsner and M. W. Godfrey, "Cloning considered harmful" considered harmful: Patterns of cloning in software," *Empirical Software Engineering*, vol. 13, no. 6, pp. 645–692, Dec. 2008.
- [12] C. K. Roy and J. R. Cordy, "A Survey on Software Clone Detection Research," *Queen's School of Computing TR*, vol. 115, p. 115, 2007.
- [13] S. Bellon, R. Koschke, G. Antoniol, J. Krinke, and E. Merlo, "Comparison and Evaluation of Clone Detection Tools," *IEEE Transactions on Software Engineering*, vol. 33, no. 9, pp. 577–591, 2007.
- [14] D. Rattan, R. Bhatia, and M. Singh, "Software clone detection: A systematic review," *Information and Software Technology*, vol. 55, no. 7, pp. 1165–1199, Jul. 2013.
- [15] E. Juergens, F. Deissenboeck, and B. Hummel, "Code Similarities Beyond Copy & Paste," in *2010 14th European Conference on Software Maintenance and Reengineering*. Madrid: IEEE, Mar. 2010, pp. 78–87.
- [16] R. Koschke, "Survey of research on software clones," in *Duplication, Redundancy, and Similarity in Software*, ser. Dagstuhl Seminar Proceedings, R. Koschke, E. Merlo, and A. Walenstein, Eds., no. 06301. Dagstuhl, Germany: Internationales Begegnungs- und Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl, Germany, 2007.
- [17] J. Svajlenko and C. K. Roy, "Evaluating clone detection tools with BigCloneBench," in *2015 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. Bremen, Germany: IEEE, Sep. 2015, pp. 131–140.
- [18] F. Farmahinifarahani, V. Saini, D. Yang, H. Sajjani, and C. V. Lopes, "On Precision of Code Clone Detection Tools," in *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, Feb. 2019, pp. 84–94.
- [19] H.-H. Wei and M. Li, "Supervised deep features for software functional clone detection by exploiting lexical and syntactical information in source code," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, ser. IJCAI'17. Melbourne, Australia: AAAI Press, Aug. 2017, pp. 3034–3040.
- [20] V. Kafer, S. Wagner, and R. Koschke, "Are there functionally similar code clones in practice?" in *2018 IEEE 12th International Workshop on Software Clones (IWSC)*. Campobasso: IEEE, Mar. 2018, pp. 2–8.
- [21] H. Thaller, L. Linsbauer, R. Ramler, and A. Egyed, "Probabilistic Software Modeling: A Data-driven Paradigm for Software Analysis," *arXiv:1912.07936 [cs]*, Dec. 2019.
- [22] H. Thaller, L. Linsbauer, and A. Egyed, "Towards Semantic Clone Detection via Probabilistic Software Modeling," in *2020 IEEE 14th International Workshop on Software Clones (IWSC)*. IEEE, 2020, pp. 64–69.
- [23] J. Krinke, "Identifying Similar Code with Program Dependence Graphs," *Proceedings Eighth Working Conference on Reverse Engineering*, pp. 301–309, 2001.
- [24] H. Thaller, R. Ramler, J. Pichler, and A. Egyed, "Exploring code clones in programmable logic controller software," in *2017 22nd IEEE International Conference on Emerging Technologies and Factory Automation (ETFA)*. Limassol: IEEE, Sep. 2017, pp. 1–8.
- [25] K. Arnold, J. Gosling, and D. Holmes, *The Java Programming Language*, 3rd ed. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2000.
- [26] H. Thaller, L. Linsbauer, and A. Egyed, "Feature Maps: A Comprehensive Software Representation for Design Pattern Detection," in *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. Hangzhou, China: IEEE, Feb. 2019, pp. 207–217.
- [27] F. J. Massey, "The Kolmogorov-Smirnov Test for Goodness of Fit," *Journal of the American Statistical Association*, vol. 46, no. 253, pp. 68–78, Mar. 1951.
- [28] G. Fink and M. Bishop, "Property-based testing: A new approach to testing for assurance," *ACM SIGSOFT Software Engineering Notes*, vol. 22, no. 4, pp. 74–80, Jul. 1997.
- [29] J. Fan, C. Zhang, and J. Zhang, "Generalized Likelihood Ratio Statistics and Wilks Phenomenon," *The Annals of Statistics*, vol. 29, no. 1, pp. 153–193, 2001.
- [30] J. R. Cordy and C. K. Roy, "The NiCad Clone Detector," in *2011 IEEE 19th International Conference on Program Comprehension*, Jun. 2011, pp. 219–220.
- [31] PMD, "PMD," PMD, 2019.
- [32] N. Göde and R. Koschke, "Incremental Clone Detection," in *2009 13th European Conference on Software Maintenance and Reengineering*, Mar. 2009, pp. 219–228.
- [33] G. Li, H. Liu, Y. Jiang, and J. Jin, "Test-Based Clone Detection: An Initial Try on Semantically Equivalent Methods," *IEEE Access*, vol. 6, pp. 77 643–77 655, 2018.
- [34] M. White, M. Tufano, C. Vendome, and D. Poshyanyk, "Deep learning code fragments for code clone detection," in *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, Sep. 2016, pp. 87–98.
- [35] G. Zhao and J. Huang, "DeepSim: Deep learning code functional similarity," in *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2018. Lake Buena Vista, FL, USA: Association for Computing Machinery, Oct. 2018, pp. 141–151.
- [36] J. D. Evans, *Straightforward Statistics for the Behavioral Sciences*. Pacific Grove: Brooks/Cole Pub. Co, 1996.
- [37] A. Walenstein, N. Jyoti, J. Li, Y. Yang, and A. Lakhotia, "Problems Creating Task-Relevant Clone Detection Reference Data," in *Proceedings - Working Conference on Reverse Engineering, WCRE*, vol. 2003-Janua, 2003, pp. 285–294.
- [38] F. Deissenboeck, L. Heinemann, B. Hummel, and S. Wagner, "Challenges of the Dynamic Detection of Functionally Similar Code Fragments," in *2012 16th European Conference on Software Maintenance and Reengineering*, Mar. 2012, pp. 299–308.
- [39] G. Li, H. Liu, Y. Jiang, and J. Jin, "Test-Based Clone Detection: An Initial Try on Semantically Equivalent Methods," *IEEE Access*, vol. 6, pp. 77 643–77 655, 2018.
- [40] S. Horwitz, "Identifying the Semantic and Textual Differences Between Two Versions of a Program," in *PLDI*, 1990.
- [41] M. Gabel, L. Jiang, and Z. Su, "Scalable Detection of Semantic Clones," in *Proceedings of the 13th International Conference on Software Engineering - ICSE '08*. ACM Press, 2008, p. 321.
- [42] R. Komondoor and S. Horwitz, "Using Slicing to Identify Duplication in Source Code," in *Proceedings of the 8th International Symposium on Static Analysis*, ser. SAS '01. London, UK, UK: Springer-Verlag, 2001, pp. 40–56.
- [43] S. Wagner, A. Abdulkhaleq, I. Bogicevic, J.-P. Ostberg, and J. Ramadani, "How are functionally similar code clones syntactically different? An empirical study and a benchmark," *PeerJ Computer Science*, vol. 2, p. e49, Mar. 2016.
- [44] L. Jiang and Z. Su, "Automatic Mining of Functionally Equivalent Code Fragments via Random Testing," in *Proceedings of the Eighteenth International Symposium on Software Testing and Analysis*, ser. ISSTA '09. New York, NY, USA: ACM, 2009, pp. 81–92.