

BAITAL: An Adaptive Weighted Sampling Approach for Improved t -wise Coverage

Eduard Baranov
Université Catholique de Louvain
Belgium
eduard.baranov@uclouvain.be

Axel Legay
Université Catholique de Louvain
Belgium
Aalborg University
Denmark
axel.legay@uclouvain.be

Kuldeep S. Meel
National University of Singapore
Singapore
meel@comp.nus.edu.sg

ABSTRACT

The rise of highly configurable complex software and its widespread usage requires design of efficient testing methodology. t -wise coverage is a leading metric to measure the quality of the testing suite and the underlying test generation engine. While uniform sampling-based test generation is widely believed to be the state of the art approach to achieve t -wise coverage in presence of constraints on the set of configurations, such a scheme often fails to achieve high t -wise coverage in presence of complex constraints. In this work, we propose a novel approach BAITAL, based on adaptive weighted sampling using literal weighted functions, to generate test sets with high t -wise coverage. We demonstrate that our approach reaches significantly higher t -wise coverage than uniform sampling. The novel usage of literal weighted sampling leaves open several interesting directions, empirical as well as theoretical, for future research.

CCS CONCEPTS

• **Software and its engineering** → **Feature interaction; Software product lines; Software testing and debugging.**

KEYWORDS

Configurable software, t -wise coverage, Weighted sampling

ACM Reference Format:

Eduard Baranov, Axel Legay, and Kuldeep S. Meel. 2020. BAITAL: An Adaptive Weighted Sampling Approach for Improved t -wise Coverage. In *Proceedings of the 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '20)*, November 8–13, 2020, Virtual Event, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3368089.3409744>

1 INTRODUCTION

The software has been one of the primary driving forces in the transformation of humanity in the past half-century; in the modern world, software touches every aspect of modern lives ranging from medical, legal, judicial to policy-making. The widespread and

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ESEC/FSE '20, November 8–13, 2020, Virtual Event, USA

© 2020 Association for Computing Machinery.
ACM ISBN 978-1-4503-7043-1/20/11...\$15.00
<https://doi.org/10.1145/3368089.3409744>

diverse usage has led to the design of highly configurable software systems operating in diverse environments. Since software failures can lead to catastrophic effects, adequate testing of configurable systems is paramount. Testing of configurable systems adds complexity on top of an already notoriously difficult problem of testing standard software.

In the context of configurable systems, every configuration refers to an assignment of values to different parameters. For the exposition, we will restrict our discussion to parameters that only take binary values; the techniques proposed in this work are general and applicable to parameters whose possible set of values form a finite set, and the benchmarks employed in our empirical study arise from such domains. The primary challenge in the testing of configurable systems arising from the observation that bugs often arise due to interactions induced by the combination of parameter values. In the combinatorial testing literature, the term feature is often used to indicate a given parameter value. One such example is an extensive study by Abal, Brabanc, and Wasowski[1] that identified 42 bugs caused by the feature combinations in the Linux kernel. Furthermore, modeling of system and environment leads to constraints over the possible set of configurations of interest.

Combinatorial testing, also known as combinatorial interaction testing (CIT), has emerged as one of the dominant paradigms for testing of configurable software wherein the focus is to employ techniques from diverse areas to generate test suites to attain high coverage. One of the widely used metrics is t -wise coverage, wherein the focus is to achieve coverage of all combinations of features of size t .

A fundamental problem in CIT is the generation of test-configurations that seeks to maximise t -wise coverage, which is measured as the fraction of feature combinations appearing in the test set out of the possible valid feature combinations. The complexity of the problem arises from the presence of constraints to capture the set of invalid configurations. The holy grail of test generation in CIT is the design of test generation methods that can handle complex constraints, scale to systems involving thousands of features, and achieve higher t -wise coverage. Since achieving *high* t -wise coverage can be infeasible for large values of t , the practitioners often focus on small values of $t \in \{1, 2, 3\}$, wherein $t = 1$ corresponds to achieving feature-wise coverage.

For long, uniform sampling has been viewed as a dominant domain-agnostic paradigm to achieve higher t -wise coverage, as demonstrated by theoretical and empirical analysis [41, 51]. As an example, the accepted solution for SPLC 2019 challenge, *Product Sampling for Product Lines: The Scalability Challenge*, was uniform

sampler, Smarch, contributed by Oh, Gazzillo, and Batory [48]. Uniform sampling seeks to sample each valid configuration with equal probability. Recent works [48, 53] present tools that are capable of performing uniform sampling and of generating a partial covering array for large feature models. Uniform sampling provides an excellent heuristic to choose samples that would cover various feature combinations; however, this approach has a limitation since feature combinations are not distributed equally among the configurations: some of them can appear in millions of configurations while others only in tens. This fact prevents uniform sampling from achieving high t -wise coverage. In [48] a feature model was shown to have only about 48% pairwise coverage achievable with uniform sampling for any reasonable number of samples, while 70% coverage required more than 10^{10} samples. It is perhaps worth highlighting strengths of the uniform sampling approach: (1) domain agnosticism, and (2) principled use of the progress in SAT solving since the problem of uniform sampling is well formulated and fundamental problem with close relationship to other problems in other problems in complexity theory [4, 8, 23]. In this context, one wonders: *whether there exists a domain-agnostic alternative to uniform sampling which can benefit from advances in formal reasoning tools?*

The key contribution of this work is an affirmative answer to the above question: We present an adaptive literal-weighted sampling approach, called BAITAL, that achieves significantly higher t -wise coverage. Our formulation builds on the recent advances in formal methods community in the development of distribution-aware sampling for the distributions specified using literal weighted functions, a rich class of distributions with a diverse set of applications. In contrast to prior sampling-based approaches that advocate sampling from a fixed distribution, BAITAL adapts the distribution based on the generated samples. Therefore, BAITAL follows a multi-round architecture wherein the initial weight function is uniform, and then the weight function is updated after each round. The update of weight function typically requires the underlying sampling algorithm to redo the entire work from scratch. To allow the reuse of the computation, we adapt the recently introduced knowledge compilation-based approach for weighted sampling [20].

We introduce several strategies for the modification of literal-weight function. Through an extensive empirical evaluation on large benchmarks, we demonstrate BAITAL achieves significantly higher pairwise coverage than uniform sampling. As an example, the accepted solution to the SPLC 2019 challenge could achieve less than 50% 2-wise coverage while BAITAL can achieve over 90% coverage with just 1000 samples. Our extensive comparison among different proposed strategies reveals surprising high coverage achieved by strategies based on limited statistics from the generated samples. The comparison is focused on pairwise coverage due to the complexity of coverage computation for higher values of t .

In summary, this work introduces a new paradigm based on adaptive weighted sampling to achieve high t -wise coverage that can benefit from advances in knowledge compilation. The high coverage obtained by BAITAL leads to several exciting directions of future work: First, we would like to develop a deeper understanding of the theoretical power of literal weight functions in the context of t -wise coverage seems a promising research area. Secondly, our work highlights challenges of measuring t -wise coverage for larger t , and the development of efficient algorithmic techniques

to estimate such coverage would be beneficial in evaluating different techniques and aid in the development of more strategies for BAITAL. Our implementation of BAITAL is publicly available at <https://doi.org/10.5281/zenodo.4028454>.

2 BACKGROUND

2.1 Boolean Formulas and Weight Function

A literal is a boolean variable or its negation. A clause is a disjunction of a set of literals. A propositional formula F in conjunctive normal form (CNF) is a conjunction of clauses. Let $Vars(F)$ be the set of variables appearing in F . The set $Vars(F)$ is called *support* of F . A *satisfying assignment* or *witness* of F , denoted by σ , is an assignment of truth values to variables in its support such that F evaluates to true. We denote the set of all witnesses of F as R_F . Let $var(l)$ denote the variable of literal l , i.e., $var(l) = var(\neg l)$ and $F|_l$ denotes the formula obtained when literal l is set to true in F .

Given a propositional formula F and a weight function $W(\cdot)$ that assigns a non-negative weight to every literal, the weight of assignment σ denoted as $W(\sigma)$ is the product of weights of all the literals appearing in σ , i.e., $W(\sigma) = \prod_{l \in \sigma} W(l)$. Without loss of generality, we assume that the weight function is normalised, i.e., $W(l) + W(\neg l) = 1$. The weight of a set of assignments \mathcal{U} is given by $W(\mathcal{U}) = \sum_{\sigma \in \mathcal{U}} W(\sigma)$. Note that, we have overloaded the definition of weight function $W(\cdot)$ to support different arguments – a literal, an assignment and a set of assignments.

2.2 t -wise Coverage

The formulation of combinatorial interaction testing (CIT) assigns a variable corresponding to every feature. Let X be the set of all the variables (corresponding to features). Furthermore, every configuration $\sigma \in 2^X$ can be represented as conjunction of the set of literals of size $|X|$, where $|\cdot|$ denotes size of a set. For example, let $X = \{x_1, x_2, x_3\}$, then $\sigma = \{x_1, x_3\}$ can be equivalently represented as $\{x_1, \neg x_2, x_3\}$

Given a configuration $\sigma \in 2^X$ represented as a set of literals, let $\text{Comb}(\sigma, t)$ represent the set of t -sized feature combinations due to σ , which is essentially the set of all subsets of literals of the size t in σ . For a set $\mathcal{U} \subseteq 2^X$, we can extend the notion of Comb and denote $\text{Comb}(\mathcal{U}, t) = \cup_{\sigma \in \mathcal{U}} \text{Comb}(\sigma, t)$. Note that while for a given σ , $|\text{Comb}(\sigma, t)| = \binom{|X|}{t}$ but this does not imply $|\text{Comb}(\mathcal{U}, t)| = |\mathcal{U}| \times \binom{|X|}{t}$ since $|\text{Comb}(\sigma_1, t) \cup \text{Comb}(\sigma_2, t)|$ is not necessarily equal to $|\text{Comb}(\sigma_1, t)| + |\text{Comb}(\sigma_2, t)|$.

t -wise coverage of a set \mathcal{U} is defined as a ratio between number of t -sized combinations due to \mathcal{U} over the total number of t -sized combinations. Feature models have constraints over a set of variables defined with a formula F , therefore for t -wise coverage only configurations that are witness of F are considered. Formally, $\text{Cov}(\mathcal{U}, t) = |\text{Comb}(\mathcal{U}, t)| / |\text{Comb}(R_F, t)|$.

Allowing only a fixed number of configurations, the coverage optimisation problem searches for a fixed-sized set of configurations that maximises the t -wise coverage. Given a formula F over X , size of feature combination t and allowed number of samples s , the problem of $\text{OptTCover}(F, t, s)$ seeks for \mathcal{U} such that:

$$\mathcal{U} = \underset{\mathcal{U} \subseteq R_F, |\mathcal{U}|=s}{\text{argmax}} |\text{Cov}(\mathcal{U}, t)|.$$

For a given F , t and s , an ideal test generator seeks to solve $\text{OptTCover}(F, t, s)$

2.3 Knowledge Compilation

We focus on subsets of Negation Normal Form (NNF) where the internal nodes are labeled with disjunction (\vee) or conjunction (\wedge) while the leaf nodes are labeled with \perp (*false*), \top (*true*), or a literal. For a node v , let $\vartheta(v)$ and $\text{Vars}(v)$ denote the formula represented by the DAG rooted at v , and the variables that label the descendants of v , respectively.

We define the well-known decomposed conjunction [12] as follows:

Definition 2.1. A conjunction node v is called a *decomposed conjunction* if its children (also known as conjuncts of v) do not share variables. Formally, let w_1, \dots, w_k be the children of \wedge -node v , then $\text{Vars}(w_i) \cap \text{Vars}(w_j) = \emptyset$ for $i \neq j$.

If each conjunction node is decomposed, we say the formula is in *Decomposable NNF* (DNNF).

Definition 2.2. A disjunction node v is called *deterministic* if each two disjuncts of v are logically contradictory. That is, if w_1, \dots, w_n are the children of \vee -node v , then $\vartheta(w_i) \wedge \vartheta(w_j) \models \text{false}$ for $i \neq j$.

If each disjunction node of a DNNF formula is deterministic, we say the formula is in *deterministic DNNF* (d-DNNF), and we can perform tractable model counting on it.

3 BAITAL: ACHIEVING HIGHER t -WISE COVERAGE

In this section, we first discuss the relationship of $\text{OptTCover}(F, t, s)$ with the classical problem of set cover. We then discuss the limitations of lifting techniques in the context of set cover to configuration testing due to presence of constraints. Inspired by the classical greedy search in the context of set cover, we design BAITAL, a novel framework that employs adaptive weighted sampling combined with recently proposed knowledge compilation-based sampling approach to achieve efficient sampling routines that achieves high coverage. The development of BAITAL leads to the need for strategies for update of weight functions; to this end, we present several strategies that seek to use different statistics from the generated sampling.

3.1 Relationship with Set Cover

As a starting point, we observe that the $\text{OptTCover}(F, t, s)$ problem can be reduced to that of the optimisation variant of the classical problem of set cover [26]; the variant is also referred to as Maximum Coverage. The reduction to set cover allows us to simply modify the classical greedy search strategy to obtain a \mathcal{U} such that $\text{Cov}(\mathcal{U}, t)$ is at least $(1 - 1/e)$ of the optimal solution \mathcal{U}^* of $\text{OptTCover}(F, t, s)$. We present the algorithm below for completeness and to discuss its simplicity and yet its impracticality in practice.

The algorithm GreedyCovertion is presented in Algorithm 1. GreedyCovertion assumes access to the subroutine MaxDistSolution that takes in input F and \mathcal{S} , and returns σ^* such that:

Algorithm 1: GreedyCovertion(F, s)

```

1  $\mathcal{S} \leftarrow \emptyset$ 
2 for  $i \leftarrow 0$  to  $s$  do
3    $\sigma \leftarrow \text{MaxDistSolution}(F, \mathcal{S})$ 
4    $\mathcal{S} \leftarrow \mathcal{S} \cup \{\sigma\}$ 
5 return  $\mathcal{S}$ 

```

$$\sigma^* = \underset{\sigma \in R_F}{\text{argmin}} |\text{Comb}(\sigma, t) \cap \text{Comb}(\mathcal{S}, t)| \quad (1)$$

The following proposition follows from classical result [43].

PROPOSITION 3.1. $\text{Cov}(\mathcal{S}, t) \geq (1 - \frac{1}{e} + o(1))\text{OptTCover}(F, t, s)$, where $(1 - \frac{1}{e} + o(1)) \approx 0.632$.

While GreedyCovertion is a simple algorithm, the roadblock lies in ensuring an efficient implementation of MaxDistSolution . Recall that the set of $|\text{Comb}(\sigma, t)| = \binom{|X|}{t} \in \Omega((n/t)^t)$, therefore, even obtaining a P^{NP} subroutine for MaxDistSolution seems a daunting challenge. We are not aware of any polynomial reduction of MaxDistSolution to polynomially many MaxSAT queries. It is worth emphasising the input size is $|F| + |\mathcal{S}|$. We leave an efficient implementation of GreedyCovertion as an open question.

3.2 Adaptive Weighted Sampling

It is worth recalling that MaxDistSolution seeks to find σ with the smallest intersection of $\text{Comb}(\sigma, t)$ and the existing $\text{Comb}(\mathcal{S}, t)$. The lack of efficient implementation of MaxDistSolution makes us wonder if one can employ a randomized sampling-based approaches to find σ that seeks to achieve the goal of MaxDistSolution .

The usage of sampling has long been explored in the context of ideal test generation. In fact, random testing is the recommended strategy [30]. The key idea is to sample solutions of F uniformly at random. The past few years is witness to several efforts to understand the scalability of uniform samplers in the context of CIT. Furthermore, in response to SPLC Challenge 2019 [50], the accepted tool Smarch generated a set of uniformly distributed samples and revealed that uniform sampling could reach only 48% pairwise coverage with 1000 samples [47]. Moreover, the coverage growth with the number of samples was prolonged, requiring a million samples to reach 50% coverage. The weak performance of uniform sampling can be attributed to an uneven distribution of feature combinations among the configurations: having 10^{14} configurations one-third of literals (and consequently all feature combinations involving them) are part of less than 10^6 configurations. The probability of covering these combinations with uniform sampling is minuscule. In this regard, we seek to find different sampling strategies that can achieve higher coverage.

The key contribution of our work is an adaptive weighted sampling-based generation of tests BAITAL. The core architecture of our framework is a multi-round process wherein each round seeks to generate samples using a weighted sampler. In contrast, the weight distribution is adjusted based on the samples generated so far. To accomplish an efficient procedure, we need to tackle three challenges:

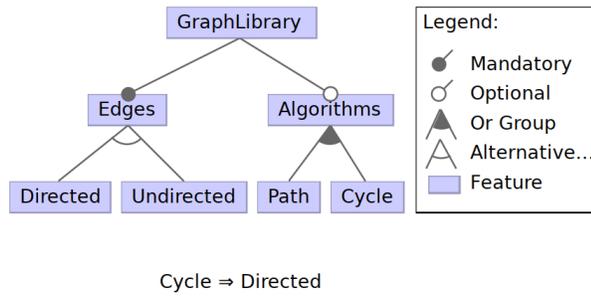


Figure 1: Feature model of graph library

Challenge 1: Representation of weights over assignment space

To represent weights over assignment space, we turn to a literal-weighted function that assigns a non-negative weight to every literal such that the weight of an assignment is the product of the weight of its literal. The choice of literal-weighted function is primarily motivated due to the observation that a wide variety of distributions arising from diverse disciplines can be represented as literal-weight function [13, 20].

Challenge 2: Dynamic update of weight function

Since the choice of literal weight function allows us to concentrate on performing the update of only the literal weights, our strategy would be to collect statistics corresponding to every literal and then assign the weights accordingly.

Challenge 3: Efficient weighted sampling techniques that can handle incremental queries

Prior techniques that seek to employ sampling techniques either fall into uniform sampling methods that solely focus on drawing uniform sampling and hope to achieve higher coverage. On the other hand, techniques that seek to induce bias in the distributions update the weights and require the underlying sampling engine to perform computations from scratch. In this work, our critical insight is to build on recent advances in knowledge compilation, a field in Artificial Intelligence that focuses on the representation of constraints in a representation language into a target language where different queries are often tractable. In particular, we seek to represent the initial set of constraints of F into a representation language \mathcal{T} such that we can perform an update of the weights and subsequent sampling in time linear in the size of \mathcal{T} without needing to perform the compilation process again. As is discussed in Section 2.3, the process of compilation is often the most expensive, and therefore, our process amortises the cost of compilation with the generation of more samples.

We illustrate the approach and show the difference from uniform sampling on an adapted version of a feature model appeared in [60]. The model shown in Figure 1 describes a product line of graph libraries. Each configuration shall always support edges that could be either D directed or U undirected and could also provide some A

algorithms like P shortest path and C deletion of cycles. The latter can only be applied on directed edges. This feature model has 6 configuration listed below, where we omit features GraphLibrary and Edges since they must be selected in all configurations:

- (1) D
- (2) $D A P$
- (3) $D A C$
- (4) $D A P C$
- (5) U
- (6) $U A P$

These configurations have 34 pairs of literals, non-selection of a feature (negation of a variable) can also be part of a feature pair. In this running example we consider that we can select 3 configurations which is not enough to obtain full pair-wise coverage. Uniform sampling selects configurations at random, all configurations have equal chances to be chosen. In the running example we assume that the first two selected configurations are (2) and (4) and we need to choose the third one. At this step configuration (3) would add only 4 new pairs, configurations (1) and (6) would add 7 pairs, while configuration (5) would add 10 pairs. At this step uniform sampling can select any of the remaining configurations with a probability 0.25.

We now first present an overview of BAITAL and then discuss each of the components in detail in the following sections.

In order to overcome the limitation of the uniform sampling, configurations with rare feature combinations shall have a higher probability of being chosen. Indeed, the precomputation of weights for all configurations is infeasible. Therefore, we are proposing to adjust the weights dynamically by examining already sampled configurations. This way, we can force the sampler to prefer configurations with uncovered feature combinations over the other ones. The algorithm runs as follows: a first few samples are generated according to the uniform distribution of configurations. After this step, the sampler pauses while the examination of the chosen configurations is done. The algorithm detects which feature combinations have been covered already and chooses the weights for the next round of sample generation accordingly. The next samples are chosen following the new weight distribution. The process repeats until the desired number of samples is generated.

For the example from Figure 1, at the first round samples (2) and (4) are generated similarly to uniform sampling. At this step the algorithm discovers that combinations involving D , A , and P have already been covered, therefore their weights shall be reduced, while the weight of U shall be raised. The modification of weights affects probability distribution during the generation of the next sample.

The overall algorithm is shown in Algorithm 2. It has three input parameters: the CNF formula F defining constraints on a feature model, the number of rounds $rounds$, and the number of samples generated at each round spr , i.e. the total number of samples to be generated is $rounds \times spr$. We first compile the formula F into a d-DNNF \mathcal{T} . At the start of each round (line 4) the algorithm generates an assignment of weights to variables that depends on the set of samples generated during the previous rounds. We discuss generateWeights in details in the next subsections. The literal-weighted sampling algorithm is called to annotate \mathcal{T} with weights (line 6) and to find spr samples chosen according to the

Algorithm 2: AdaptiveWeightedSampling(F , $rounds$, spr)

```

1  $\mathcal{S} \leftarrow \emptyset$ 
2  $vars \leftarrow$  Set of variables in  $F$ 
3  $\mathcal{T} \leftarrow Compile(F)$ 
4 for  $roundNb \leftarrow 0$  to  $rounds$  do
5    $weights \leftarrow generateWeights(vars, \mathcal{S})$ 
6   Annotate( $\mathcal{T}$ ,  $weights$ )
7    $newSamples \leftarrow Sample(\mathcal{T}, spr)$ 
8    $\mathcal{S} \leftarrow \mathcal{S} \cup newSamples$ 
9 return  $\mathcal{S}$ 

```

distribution imposed by generated weights (line 7). For the first round, the weights are always equal to 0.5, corresponding to the uniform distribution. The new samples are added to the result (line 8), and the procedure repeats $rounds$ times.

3.3 Weighted Sampling via Compilation

In [53] it was observed that advances in knowledge compilation can be used to design efficient uniform samplers and designed an efficient uniform sampler KUS. In [20] KUS was extended to design a weighted sampler WAPS. Given an input formula F and the desired number of samples s , WAPS uses a three staged process: **Compilation**, **Annotation**, and **Sampling**.

Compilation Given a formula F , a state of the art compiler is employed to obtain the equivalent d-DNNF \mathcal{T} .

Annotation A bottom-up subroutine is called to annotate every node v of \mathcal{T} with its corresponding weight, i.e., the label for the node v is $W(\vartheta(v))$. It is worth highlighting that for a formula in d-DNNF, the labeling of every node with its weight can be accomplished in time linear in the size of \mathcal{T} .

Sampling To perform sampling, a top-down subroutine is followed wherein for every \wedge node, once the samples from all the children are generated, we randomly shuffle each of the lists and then conjoin the list of samples. In the case of \vee node, we first determine, via the Binomial coefficient, the number of samples each of the nodes should generate, and then we take the union of all the lists.

It is known that there exist polynomially sized formulas whose d-DNNF representations are exponential [12]. Furthermore, the compilation process may require a significantly larger time than the size of resulting \mathcal{T} . For example, for every formula F that is valid, there exists a polynomially sized d-DNNF \mathcal{T} , but the existence of PTIME compilation process would imply P=Co-NP. Fortunately, our adaptive sampling strategy only modifies the weight function and does not alter F . We make the following key observation about **Annotation** and **Sampling** in WAPS that strongly supports our choice of knowledge compilation-based approach for BAITAL.

PROPOSITION 3.2. *Annotation and Sampling can be accomplished in time linear in the size of \mathcal{T}*

PROOF. The proof follows trivially from the algorithmic description of **Annotation** and **Sampling** in [20, 53]. Essentially, the key observation is that **Annotation** is performed in a bottom-up process while visiting each node exactly once, while **Sampling** is

performed in a top-down manner while visiting each node exactly once. \square

3.4 Round Weights Generation

In literal-weighted sampling, configurations are chosen according to the distribution imposed by a weight function W . As shown in Section 2, W is defined by assigning a (normalised) weight to each literal.

For the first round the weight function is constant, $W(l) = 0.5$ for any literal l which corresponds to the uniform distribution. For the following rounds the definition of W is based on the knowledge about constraints F , and on a set of already generated samples \mathcal{S} . Let $g(l, t, \mathcal{S})$ be a function that represents the knowledge about the t -sized feature combinations involving a literal l in the set of samples \mathcal{S} . Let $h(l, t, F)$ be a function that represents the knowledge about the t -sized feature combinations involving a literal l by configurations in R_F . We define the weight function as follows:

$$W(l, t, \mathcal{S}, F) = f(g(l, t, \mathcal{S}), h(l, t, F), g(\neg l, t, \mathcal{S}), h(\neg l, t, F)), \quad (2)$$

where F is a propositional formula, l is a literal, t is a size of feature combinations, \mathcal{S} is a set of already generated samples, and f is a function outputting a value in the interval $[0, 1]$. Since the weights of a literal and its negation must be related, the weight function is dependent on the feature combinations of the literal negation as well. Each round of BAITAL starts with the computation of W (line 5 of Algorithm 2) which is further used for sample generation.

There exist multiple ways to define functions f , g , and h resulting in different weight functions and, consequently, in different t -wise coverage. We have considered several options to define these functions, which we call strategies. In the remaining, we assume that F and t are fixed and omit them. Notice that $h(l)$ depends only on fixed F and t and, therefore, is not changing between the rounds and can be computed only once. For illustration we use the running example from Figure 1 where we have selected configurations (2) and (4) and compute the weights for the next round.

Algorithm 3: generateWeights($vars$, \mathcal{S}) // Strategy 1

```

1 for  $l$  in literals do
2    $g[l] \leftarrow$  number of distinct  $t$ -sized feature combinations
   with  $l$  in  $\mathcal{S}$ 
3    $h[l] \leftarrow$  number of distinct  $t$ -sized feature combinations
   with  $l$  in  $R_F$ 
   /*  $h$  is computed in the first round only */
4 for  $l$  in literals do
5    $diff[l] \leftarrow g[l] / h[l] - g[\neg l] / h[\neg l]$ 
6    $weights[l] \leftarrow max(0.01,$ 
    $0.5 * (1 - sign(diff) * sqrt(abs(diff))))$ 
7 return weights

```

Strategy 1. The first option is to define g and h as the number of distinct feature combinations involving each literal in a set of samples and in R_F , respectively. For the computation of h it is not necessary to enumerate all configurations in R_F , it can be done by checking, for each feature combination $\{l_1, \dots, l_t\}$, satisfiability of

$F \wedge l_1 \wedge \dots \wedge l_t$. Indeed, since the SAT solver can provide an assignment in case the formula is satisfiable, during the computation of h , the full covering array is implicitly built. However, for 1000 variables, such array would contain millions of samples, thus being impractical.

The heuristic behind this strategy is the following: if most of feature combinations involving a literal are already covered but it is not the case for the negation of the literal (or vice versa), then the sampler should prefer configurations involving a negation of a literal in the next round. We use the difference between ratios of covered feature combinations for literal and its negation as a defining value for the weight. Therefore, function f is computing this difference and transforms it into the interval $[0, 1]$. Algorithm 3 illustrates the weight generation for strategy 1 with a pseudo-code. First, g and h are computed over lines 2 and 3. In the following step, f is computed over lines 5 and 6. The choice of sqrt-based transformation function line 6 is done after the empirical evaluation of several options. We use max ensures that the resulted weight is not equal to 0.

Strategy 1 involves a lot of computations for functions g and h . In particular, h checks satisfiability of $2^t * \binom{Vars(F)}{t}$ formulas. Secondly, counting the number of distinct combinations in the set of samples has a high time and memory costs for large real-world models involving thousands of variables. Therefore, we considered several modifications of the weight function that would require fewer resources for computation and evaluated their effect on t -wise coverage.

In the running example, we compute that in samples (2) and (4) D has pairs with $A, P, C, -C, -U$ covering 5 combinations, i.e. $g[D] = 5$. During precomputation we learn that D and U do not appear together in any configurations, therefore $h[D] = 7$. At the same time $-D$ has not been covered, so $g[-D] = 0$. Applying the Algorithm 3, new weight of D is 0.08. Similar computations provide the following results for other features: $weights[A] = 0.1$, $weights[P] = 0.08$, $weights[C] = 0.37$, $weights[U] = 0.92$. Note that due to normalisation weights for negations of variables can be computed from $weights[l] + weights[-l] = 1$. Applying the extension of weight function to configurations, we can compute the probabilities of configurations to be selected. These weights correspond to 0.83 probability to choose configuration with the best additional coverage (5) as the next sample in comparison to 0.25 probability at uniform sampling. The worst next sample (3) can be chosen with just 0.004 probability.

Algorithm 4: generateWeights($vars, S$) // Strategy 2

```

1 for  $l$  in literals do
2    $g[l] \leftarrow$  number of distinct  $t$ -sized feature combinations
   with  $l$  in  $S$ 
3    $h[l] \leftarrow 2^{t-1} * \binom{Vars(F)-1}{t-1}$ 
4 for  $l$  in literals do
5    $diff[l] \leftarrow g[l] / h[l] - g[-l] / h[-l]$ 
6    $weights[l] \leftarrow max(0.01,$ 
    $0.5 * (1 - sign(diff) * sqrt(abs(diff))))$ 
7 return weights
```

Strategy 2. In this strategy we attempt to reduce the computation cost by simplifying h . Without checking the existence of configurations in R_F involving each feature combination, h can be overapproximated by the number of t -sized feature combinations in $2^{Vars(F)}$ involving a literal. The generateWeights function for the **Strategy 2** is shown in Algorithm 4. The only difference from Algorithm 3 is on line 3, since f and g are defined as in **Strategy 1**.

In the running example, we do not precompute that $h[D] = 7$ in this strategy, but we get the value 8 as overapproximation. This changes the resulted weight distribution, the Algorithm 4 returns the following weights: $weights[D] = weights[A] = weights[P] = 0.1$, $weights[C] = 0.5$, $weights[U] = 0.9$. These weights correspond to 0.81 probability to choose the best next sample (5) (0.83 with **Strategy 1**) and 0.01 to choose the worst (0.004 with **Strategy 1**).

Algorithm 5: generateWeights($vars, S$) // Strategy 3

```

1 for  $l$  in literals do
2    $g[l] \leftarrow |\sigma \in S \mid l \in \sigma|$ 
3    $h[l] \leftarrow |\sigma \in R_F \mid l \in \sigma|$ 
4 for  $l$  in literals do
5    $diff[l] \leftarrow g[l] / (g[l] + g[-l]) -$ 
    $ln(h[l]) / (ln(h[l]) + ln(h[-l]))$ 
6    $weights[l] \leftarrow max(0.01,$ 
    $0.5 * (1 - sign(diff) * sqrt(abs(diff))))$ 
7 return weights
```

Strategy 3. Both **Strategies 1** and **2** are trying to optimise the coverage for a particular size t of feature combinations. Indeed, the computation of g and h depends on t , and for the same set of samples, the weights generated for $t = 2$ and $t = 3$ would be different. Nevertheless, one could expect that a set of samples with high pairwise coverage would also have high 3-wise coverage. This hypothesis allows us to gradually simplify the computation of g .

A very simple and easily computable measure of literal participation in a set of samples is the number of samples the literal is present. It allows finding the number of feature combinations the literal is involved in, though not the number of distinct combinations. Similarly, h in this strategy is defined as the total number of configurations involving the literal. h can be computed by calling #SAT $Vars(F) + 1$ times (any configuration involves either a literal or its negation but not both). Since the values of h and g can differ by several orders of magnitude, we changed f : it is comparing the ratio of samples involving the literal in the sample set and the ratio of configurations involving a literal on a logarithmic scale. The weight generation function is shown in Algorithm 5.

In the running example, we see that D appears in both samples (2) and (4). For the computation of h , we precompute that D appears in 4 configurations and $-D$ in the 2 remaining ones. Applying the Algorithm 5, new weight of D is 0.22. Similar computations for other features yield: $weights[A] = 0.22$, $weights[P] = 0.15$, $weights[C] = 0.3$, $weights[U] = 0.79$. These weights correspond to 0.88 probability to choose configuration (5) as the next sample and 0.008 to choose configuration (3).

Algorithm 6: generateWeights($vars, S$) // Strategy 4

```

1 for  $l$  in literals do
2    $g[l] \leftarrow |\sigma \in S \mid l \in \sigma|$ 
3 for  $l$  in literals do
4    $weights[l] \leftarrow \max(0.01,$ 
       $0.5 * (1 - (g[l] - g[\neg l]) / (g[l] + g[\neg l])))$ 
5 return weights

```

Strategy 4. In this strategy we modify **Strategy 3** by fixing h constant. The hypothesis of this strategy considers that since values of g and h differ by several orders of magnitude, their direct or indirect comparison might not be useful. Therefore, in this strategy, we simply compute g and compare the appearances of each literal with its negation. The function is shown in Algorithm 6.

In the running example, we see that D appears in both samples (2) and (4). Applying the Algorithm 6, new weight of Directed is 0.01. New weights for other features are: $weights[A] = weights[P] = 0.01$, $weights[C] = 0.5$, $weights[U] = 0.99$. These weights correspond to 0.9998 probability to choose configuration (5) as the next sample and only 1.0×10^{-6} probability to choose configuration (3).

4 EXPERIMENTS

To analyse the efficiency of our approach we have implemented a prototype of BAITAL in Python¹. The sampling process is done by the literal-weighted sampler WAPS[20]. To evaluate our approach we designed a set of experiments helping to answer the following research questions.

Our main goal is to push forward the t -wise coverage of the state-of-the-art approach for large SPLs, therefore the first two questions naturally arise.

RQ1 Can BAITAL be used to generate partial covering arrays for large SPLs?

RQ2 Can BAITAL achieve higher t -wise coverage than uniform sampling?

In addition, we are also interested in the evaluation of effects of parameters in our approach on the efficiency and performance.

RQ3 How often shall weight function be regenerated in order to obtain best coverage/performance?

RQ4 What is the impact of different strategies on coverage and performance?

4.1 Benchmarks

We took a large number of publicly available feature models from real-world configurable systems that were used before in evaluation of uniform sampling tools. In particular, we took the all non-synthetic benchmarks appearing in [29, 34, 50, 51] with exception of few largest feature models on which both uniform sampling and our approach run out of time or memory. Also we excluded small models that have less than 500 variables. The resulted benchmark set consists of 124 feature models². The size of feature models varies between 565 and 11254 variables, between 1164 and 62183 clauses and between 9.7×10^{13} and 7.7×10^{417} .

¹<https://github.com/meelgroup/baital>

²<https://doi.org/10.5281/zenodo.4022395>

For the last two research questions we performed a more detailed exploration on 3 benchmarks from our set. These models are: "ecos-icse11" with 1244 variables and 3146 clauses from [51], "embtoolkit" with 2128 variables and 15483 clauses from [29] (transformed into CNF formulas with FeatureIDE 3.6.0³) and "FinancialServices01" version "2018-05-09" ("financial" in the rest of the paper) with 771 variables and 7241 clauses from [50]. The latter one is particularly interesting since uniform sampling cannot achieve high coverage even with 10^7 samples (below 70% for 1-wise coverage and below 50% for 2-wise coverage) [48].

4.2 Experiments, Settings, and Competitors

We designed two experiments for the evaluation of BAITAL. The first experiment is used to answer **RQ1** and **RQ2**. In this experiment, we generated 1000 samples with both uniform sampling and BAITAL for each feature model from our benchmark set. In BAITAL the number of rounds was set to 10, each round 100 samples have been generated. We used **Strategies 2** and **4** since they do not require additional precomputations of h on each benchmark. Indeed, the computation of all feasible pairwise combinations for uClinux-config feature model would require 250 millions of SAT solver calls. For the same reason, we do not compare the t -wise coverage on these benchmarks but the number of distinct feature combinations in the sample sets, i.e. $|\text{Comb}(S, t)|$. In addition, we limited the exploration to the pairwise coverage as the computation of Comb for higher-wise coverage were exceeding available RAM on all but smallest benchmarks.

In the preparation of experiment we tried different tools for uniform sampling, namely QuickSampler [16], Smarch [47], KUS[53] and WAPS with weights corresponding to the uniform sampling. Quicksampler despite being fast checks the validity of samples at the very last step, therefore the number of valid samples is usually lower than the requested one. In many cases it generated 1000 samples none of which were valid. Also it was shown in [51] that the results of Quicksampler are often far from uniform. The remaining 3 tools were able to generate uniformly distributed sets of samples and, unsurprisingly, the resulted t -wise coverage was almost identical. Comparing their execution time, KUS was the fastest, closely followed by WAPS, while Smarch was considerably slower by three orders of magnitude on some of the benchmarks. Considering a minor difference between KUS and WAPS, we decided to use the latter one since it would allow us to have better evaluation of the additional complexity of our approach.

RQ3 and **RQ4** are related to the evaluation of BAITAL parameters. In the second experiment we used "ecos-icse11", "embtoolkit" and "financial" benchmarks. Within this experiment we generated 3000 samples with each strategy from Section 3.4 and values for the number of samples generated at each round spr were taken from the set {25, 50, 100, 200, 300}. In this experiment we checked pairwise coverage and execution time. The results of uniform sampling were used as a baseline.

Due to the probabilistic nature of both uniform sampling and BAITAL, all experiments have been run 5 times without fixed random seeds and the reported results show the mean values among 5 runs. Nevertheless, for 108 out of 124 benchmarks, the difference

³<https://github.com/FeatureIDE/FeatureIDE/releases/tag/v3.6.0>

Table 1: Number of feature combinations in 1000 samples

Benchmark	NVars	Uniform	Strategy 2	Strategy 4
busybox_1_28_0	998	1936827	1952827	1955269
dreamcast	1252	2308471	2650081	2842723
ecos-icse11	1244	2299262	2674216	2842662
freebsd-icse11	1396	3564507	3712877	3720593
integrator_arm9	1267	2337336	2737317	2926949
mpc50	1213	2169357	2508473	2670210
ocelot	1266	2339932	2721531	2921261
pc_i82544	1259	2291769	2741419	2887288
refdt334	1263	2330423	2765357	2941096
XSEngine	1260	2342488	2718346	2901167

between the highest and the lowest values of Comb for the same input parameters was below 1%; the maximum difference was 1.8% on financial benchmark.

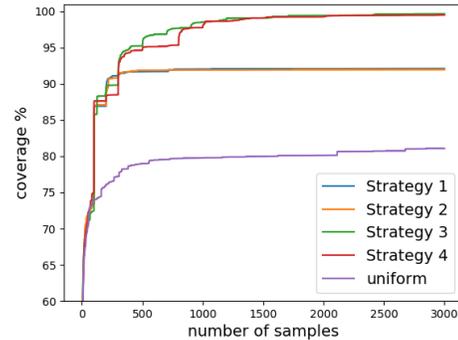
All experiments were conducted on a laptop with Intel(R) Core i7-8650U CPU, RAM limited to 12GB, running Ubuntu 18.04.3 LTS.

4.3 Comparison with Uniform Sampling (RQ1, RQ2)

The first experiment compared the number of feature combinations in sets of samples generated by uniform sampling and by BAITAL with **Strategies 2 and 4**. Table 1 presents the excerpt of the results for 10 feature models showing the number of distinct feature pairs in 1000 generated samples⁴. Among 124 benchmarks, on 116 of them the **Strategy 2** showed the increase between 14% and 22% of the number of covered feature pairs in comparison with uniform sampling and **Strategy 4** showed increase between 21% and 29%. One benchmark showed increase by 37% with **Strategy 2** and 49% with **Strategy 4**. On 5 benchmarks smaller increase was encountered and 2 remaining benchmarks returned identical results. Detailed exploration of the latter 7 benchmarks showed that on 5 of them uniform sampling reached coverage of at least 93% (for the remaining 2 we were not able to compute $|\text{Comb}(R_F, t)|$) that explains smaller improvement of BAITAL on these benchmarks. For example, busybox_1_28_0 has 98.5% coverage with uniform sampling, and BAITAL can still improve the result to 99.5%, though the difference in the number of feature combinations is only 1%.

To evaluate the presence of statistical difference between uniform sampling and BAITAL we used Mann-Whitney U-test with the H_0 : uniform sampling distribution and BAITAL sampling distribution are equal. We have a sample size 5 for both approaches, $\alpha = 0.05$. The critical value to reject H_0 is 2, i.e. the hypothesis cannot be rejected if on 3 pairwise comparisons uniform sampling would outperform BAITAL. On two benchmarks with identical results the U-test confirmed the hypothesis, while on all other benchmarks the H_0 have been rejected with $U = 0$ and $pvalue = 0.004$.

The execution time of BAITAL depends on 3 major factors: sampling time of each round, computation of weight function for the following round and the number of rounds. The first factor is the

**Figure 2: Number of feature combinations for ecos-icse11 benchmark**

execution time of WAPS. As explained in Section 3.3, WAPS execution consists of 3 steps: Compilation, Annotation, and Sampling. Since the formula is not modified during the sampling process, Compilation is done only in the first round. On the benchmarks this step takes several seconds with a highest value of 26 seconds and median 6.5 seconds.

The computation of the weight function for the next round depends on the strategy and the number of variables in the formula. For **Strategies 3 and 4**, as they just compute the number of appearances of a literal in samples, it takes less than one second even on the largest benchmark with 11254 variables. **Strategies 1 and 2**, on the other hand, need to find distinct feature combinations and their execution time also depends on the size of feature combinations. For pairwise sampling, it takes 1 second for 771 variables, 16 seconds for 1396 variables and up to 250 seconds for the largest benchmark with 11254 variables. Note that the set of distinct feature combinations do not need to be recomputed from scratch at each round, the results of first rounds can be reused in the latter ones. Therefore, the time does not increase with growth of the round number.

For the benchmarks, 102 have finished within 1 hour for both strategies, another 16 finished within 2 hours. The remaining 5 benchmarks with the longest Annotation process took several hours. The difference between **Strategy 2** and **Strategy 4** was less than 15 minutes on 110 benchmarks and the maximum was 35 minutes for the benchmark with 11254 variables.

As a result of this experiment we can conclude that BAITAL can generate partial covering arrays at reasonable time for large feature models. The resulted coverage of feature combinations is higher than with uniform sampling on all but 2 benchmarks where equal results were recorded.

4.4 Comparison of Weight Generation Strategies (RQ3, RQ4)

Figures 2, 3, 4 show the coverage on three benchmarks achieved with each strategy. The results of uniform sampling are shown for comparison. For all strategies, the weight function has been changed after every 100 samples and the coverage was computed after each generated sample. The first 100 samples are generated uniformly

⁴Full version of Table 1 can be found at <https://doi.org/10.5281/zenodo.4022395>

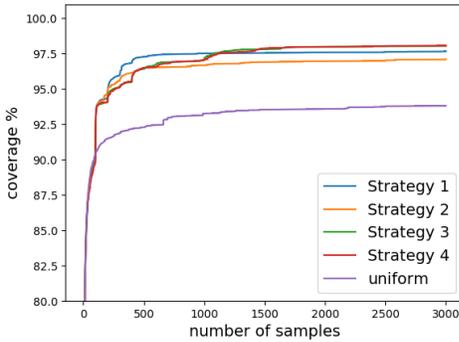


Figure 3: Number of feature combinations for embtoolkit benchmark

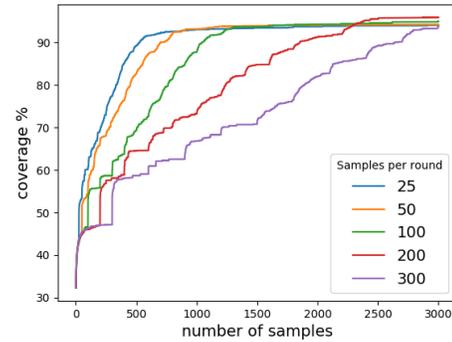


Figure 5: Number of feature combinations for financial benchmark with Strategy 1

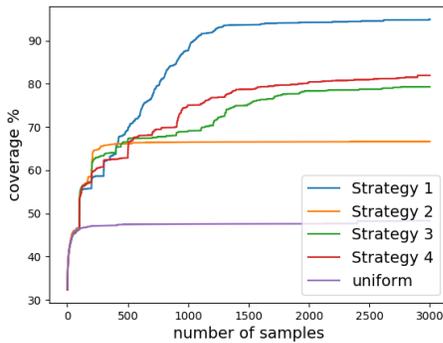


Figure 4: Number of feature combinations for financial benchmark

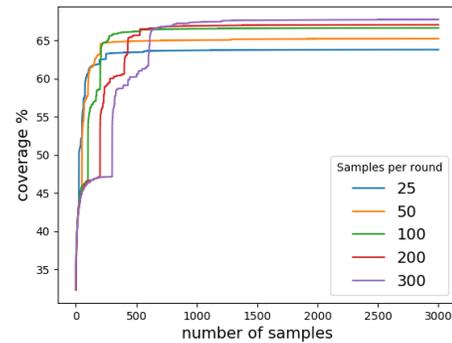


Figure 6: Number of feature combinations for financial benchmark with strategy 2

in all strategies, resulting in almost identical coverage. Starting from the second round, next sets of samples have been chosen with respect to the new weights generated by BAITAL strategies. The considerable boost over the uniform sampling can be noticed: new weights are helping to choose samples with feature combinations that didn't appear in samples before. In all three benchmarks, 200 samples generated with BAITAL have better coverage than 3000 samples chosen uniformly.

The difference between strategies becomes noticeable after second change of weights. In the first two benchmarks, the coverage is quite high even for the uniform sampling and results of different strategies are close. However, in the financial benchmark **Strategy 1** clearly outperforms other strategies, while **Strategy 2** showed the worst result. In this benchmark almost 25% of feature combinations are infeasible and their distribution among literals is far from uniform. This results in bad approximation of function $h(l, t)$ and, consequently, low performance of **Strategy 2**. **Strategies 3** and **4** showed very close results in all 3 benchmarks, therefore we can conclude that extra precomputations in **Strategy 3** do not improve to the resulted coverage.

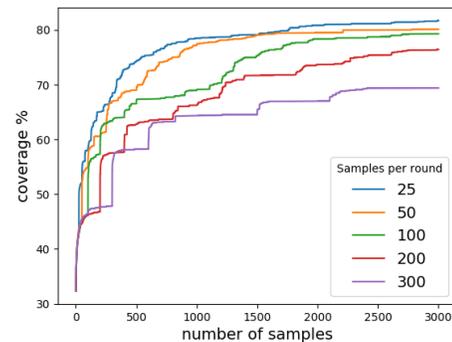


Figure 7: Number of feature combinations for financial benchmark with strategy 4

The evaluation of the effect of the weights update frequency is shown in Figures 5, 6, 7. We used the financial benchmark and we omit the plot for **Strategy 3** since it is similar to Figure 7. These plots show that higher frequencies allows to obtain high coverage

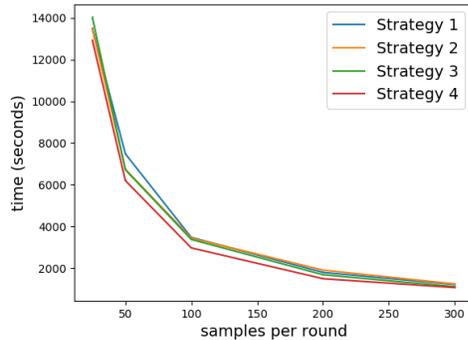


Figure 8: Execution time for generation of 3000 samples for ecos-icse11 benchmark. Uniform sampling took 107 seconds.

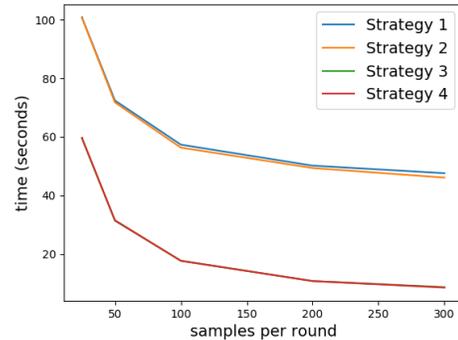


Figure 10: Execution time for generation of 3000 samples for financial benchmark. Uniform sampling took 4 seconds.

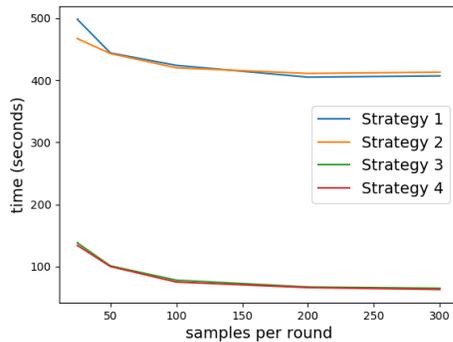


Figure 9: Execution time for generation of 3000 samples for embtoolkit benchmark. Uniform sampling took 61 seconds.

with fewer samples. However, by raising the number of generated samples, lower frequencies of weight updates can reach the same coverage or even slightly outperform it. This behaviour is a consequence of our weight generation. The difference between the weight functions of two consecutive rounds is small if during the former round only few new feature combinations have been covered. Therefore, for the first several rounds, when a lot of new feature combinations are found, high frequency of weight updates is beneficial to reach uncovered areas faster, while for the latter rounds, there would not be enough new combinations to noticeably modify weights distribution.

The evaluation of execution time on the benchmarks is shown on Figures 8, 9 and 10. The precomputations of function h are not included in the total time. Strategies using the same function g showed close results. On financial and embtoolkit benchmarks **Strategies 1 and 2** are considerably slower than the other ones and the time difference depends on the number of variables. For ecos-icse11 benchmark the annotation of d-DNNF took 95% of round time and the different strategies to generate weights had almost no effect on the total time. Regarding the frequency of weight function

generation, clearly, extra d-DNNF dag reannotations has negative impact on execution time. For feature models with relatively simple d-DNNF dag, such as embtoolkit, there is no strong effect of this parameter while in ecos-icse11 with complex d-DNNF the dependency of execution time on the number of rounds is close to linear.

As a result of this experiment, we can conclude that **Strategy 1** can generate the best partial covering array while being the slowest one. **Strategy 4** is the fastest one and provide a good coverage outperforming **Strategies 2 and 3**. For the frequency of weights change, to generate a small number of samples with high coverage, high frequency of weight function generation is preferable, however for larger sample sets it might be better to reduce it, as similar coverage could be obtained with shorter execution time.

4.5 Threats to Validity

Construct Validity. Many prior work use combinatorial interaction testing for feature models and t -wise coverage as qualitative metric for partial covering arrays for the cases where the model is too large and complex to compute a full covering array. Our approach is following these works by building partial covering arrays with better t -wise coverage.

Internal Validity. Due to the probabilistic nature of our approach, several runs may not yield identical results. Indeed, the choice of samples during the first round affects the weights for the following round. To mitigate the effect of this behaviour on the experiment results, we run each benchmark multiple times and report the mean value. In addition, we examined the result of different runs and noticed, that the difference between the best and the worst result was below 1% for 108 out of 124 of benchmarks and maximum value was 1.8%. We also noted that even the worst result was always better than the uniform sampling except two benchmarks where equality was obtained. For the baseline we used uniform sampling which is a state-of-the-art approach and we have tried several tools for uniform sampling that yield similar results.

External Validity. To mitigate the threat of non-generalisability of our study we have used a large number of feature models. These models cover a wide range in the number of features and constraints. These benchmarks were used before in several prior studies [29, 34, 48, 51].

5 RELATED WORK

Combinatorial interaction testing was initially proposed in 1980s in [38, 58] and has since been extensively studied. Covering arrays formally defined in [56] are matrices with rows representing samples and every possible t -sized combination of variables appears at least in one column. Many works on the topic were surveyed in [44] and a more recently in [59].

The goal in classical CIT sampling is to build a covering array and minimise its size. Multiple approaches have been proposed through years for building covering arrays. The greedy algorithm is one of the most popular ideas starting from the empty set of samples and adding new ones until the full coverage is achieved. There is a number of tools available including [2, 11, 24, 39, 57, 62]. Another type of greedy algorithm initially build a set of samples with full coverage for the first n parameters and then add an extra parameter at each interaction. IPO algorithm [33] with several extensions [31, 32] is following this idea as well other some works [9, 61]. Another group of approaches, sometimes called heuristical, start from some set of samples that do not provide full coverage and try to modify it in order to obtain the full coverage. Examples of such algorithms are [49], tabu search [19, 45] and genetic algorithm [40]. Lin et al. [35] attempted to combine heuristic and greedy approaches switching between them with a predefined probability, which allows a better exploration of configuration space and potentially smaller test suite. The requirement of explicit access to the entire set of configurations severely limits the above set of works in their ability to handle a large set of features.

Covering arrays may grow large for complex configurable systems and exceed the testing budget. An alternative to the classical CIT is building a fixed-size set of configurations. Several metrics have been proposed to measure combinatorial coverage of such sets in the NIST report [30]. t -wise coverage has been defined in the report as (total) variable-value configuration coverage. One of the approaches to build a fixed-size test-set is based on selecting solutions for a given set of constraints. Among constraint-based approaches, one set of techniques often seek to generate a random configuration and then check if the generated configuration satisfies the constraints [10, 15]. Such approaches can handle scenarios where the fraction of valid configurations over the set of all configurations is high but fail to handle cases where the density of valid configuration is low as typically observed in complex systems. It is worth noting that a low density of valid configurations does not imply a low number of valid configurations. For example, for 50 binary features, even if the set of valid configurations is 2^{40} , the density of such a set is still 2^{-10} .

Another set of constraint-based approach seeks to rely on employing uniform samplers constructed on top of recent advances in combinatorial solving such as SAT solving and knowledge compilation. Oh et al. [46] employed Binary Decision Diagrams for encoding the configurations of feature models. The tool SMARCH [47] uses #SAT solver for uniform sample generation. Lopez-Herrejon et al. [36] evaluated the dependency between t -wise coverage and the number of samples, building a Pareto front. Two uniform sampling tools QuickSampler [16] and Unigen2 [6] have been evaluated in [51].

Other techniques to build sample sets and their flavors have been surveyed in [3]. Several other approaches to choose samples for analysis of feature models, including t -wise sampling, one-disabled [1] and statement-coverage [57] were compared in [41].

Another application of the sampling of configurable systems is performance prediction [52, 54, 55] as different features and their interactions can affect the performance. In recent work, Kaltenecker et al. [25] proposed a new metric and the corresponding sampling method tailored for performance prediction.

The problem of weighted sampling has witnessed sustained interest from theoreticians and practitioners alike for the past three decades owing to its widespread usage. Of several proposed techniques, Monte Carlo Markov Chain (MCMC)-based methods [22, 37] enjoy widespread adoption owing to their simplicity but the heuristics employed to act as a proxy of mixing of the underlying Markov Chains, however, invalidate distributional guarantees [27]. Approaches based on interval propagation and belief networks [14, 18, 21], also often lead to scalable techniques, but their distributions can often be far from the desired distribution [28]. While hashing-based techniques have achieved significant progress in the context of uniform sampling [7, 8, 42], their scalability for arbitrary distributions is severely limited [5, 6, 17, 42]. In this context, we advocate the usage of knowledge compilation-based techniques, which not only achieve significant scalability but supports efficient adaptive sampling.

6 CONCLUSION

Design of test generation techniques to achieve higher t -wise coverage is a fundamental challenge in the context of testing of configurable systems. In this work, we propose a BAITAL approach for construction of test sample sets that dynamically modifies the probability distribution of configuration to be chosen in order to improve the t -wise coverage of the test set. We showed on a large set of benchmarks used in prior works that our approach generates achieves higher t -wise coverage than commonly used uniform sampling. We explored several ways to update the weight function and their effect on the resulted coverage.

Since our approach generates samples for combinatorial interaction testing, an interesting direction of future work would be to study the impact of improved t -wise coverage on fault detection.

ACKNOWLEDGEMENTS

This work was supported in part by EU H2020 project Serums (826278-SERUMS-H2020-SC1-FA-DTS-2018-2020) and by National Research Foundation Singapore under its NRF Fellowship Programme [NRF-NRFFAI1-2019-0004] and AI Singapore Programme [AISG-RP-2018-005], and NUS ODPRT Grant [R-252-000-685-13]. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore.

REFERENCES

- [1] Iago Abal, Claus Brabrand, and Andrzej Wasowski. 2014. 42 variability bugs in the linux kernel: a qualitative analysis. In *Proceedings of the 29th ACM/IEEE international conference on Automated software engineering*. 421–432. <https://doi.org/10.1145/2642937.2642990>
- [2] Mustafa Al-Hajjaji, Sebastian Krieter, Thomas Thüm, Malte Lochau, and Gunter Saake. 2016. Incling: efficient product-line testing using incremental pairwise

- sampling. *ACM SIGPLAN Notices* 52, 3 (2016), 144–155. <https://doi.org/10.1007/BF01588971>
- [3] Saswat Anand, Edmund K Burke, Tsong Yueh Chen, John Clark, Myra B Cohen, Wolfgang Grieskamp, Mark Harman, Mary Jean Harrold, Phil McMinn, Antonia Bertolino, et al. 2013. An orchestrated survey of methodologies for automated software test case generation. *Journal of Systems and Software* 86, 8 (2013), 1978–2001. <https://doi.org/10.1016/j.jss.2013.02.061>
 - [4] Mihir Bellare, Oded Goldreich, and Erez Petrank. 2000. Uniform Generation of NP-witnesses using an NP-oracle. *Information and Computation* 163, 2 (2000), 510–526. <https://doi.org/10.1006/inco.2000.2885>
 - [5] Supratik Chakraborty, Daniel J. Fremont, Kuldeep S. Meel, Sanjit A. Seshia, and Moshe Y. Vardi. 2014. Distribution-Aware Sampling and Weighted Model Counting for SAT. In *Proc. of AAAI*. 1722–1730.
 - [6] Supratik Chakraborty, Daniel J. Fremont, Kuldeep S. Meel, Sanjit A. Seshia, and Moshe Y. Vardi. 2015. On Parallel Scalable Uniform SAT Witness Generation. In *Proc. of TACAS*. 304–319. https://doi.org/10.1007/978-3-662-46681-0_25
 - [7] Supratik Chakraborty, Kuldeep S. Meel, and Moshe Y. Vardi. 2013. A Scalable and Nearly Uniform Generator of SAT Witnesses. In *Proc. of CAV*. 608–623. https://doi.org/10.1007/978-3-642-39799-8_40
 - [8] Supratik Chakraborty, Kuldeep S. Meel, and Moshe Y. Vardi. 2014. Balancing Scalability and Uniformity in SAT Witness Generator. In *Proc. of DAC*. 1–6. <https://doi.org/10.1145/2593069.2593097>
 - [9] Nie Changhai, Xu Baowen, Shi Liang, and Wang Ziyuan. 2006. A new heuristic for test suite generation for pair-wise testing. *SEKE 2006* 1 (2006), 517.
 - [10] Tsong Yueh Chen, Hing Leung, and IK Mak. 2004. Adaptive random testing. In *Annual Asian Computing Science Conference*. Springer, 320–329. https://doi.org/10.1007/978-3-540-30502-6_23
 - [11] Anastasia Cmyrev and Ralf Reissing. 2014. Efficient and effective testing of automotive software product lines. *King Mongkuts University of Technology North Bangkok International Journal of Applied Science and Technology* 7, 2 (2014), 53–57. <https://doi.org/10.14416/j.ijast.2014.05.001>
 - [12] Adnan Darwiche and Pierre Marquis. 2002. A knowledge compilation map. *Journal of Artificial Intelligence Research* 17 (2002), 229–264. <https://doi.org/10.1613/jair.989>
 - [13] Alexis de Colnet and Kuldeep S Meel. 2019. Dual Hashing-Based Algorithms for Discrete Integration. In *International Conference on Principles and Practice of Constraint Programming*. Springer, 161–176. https://doi.org/10.1007/978-3-030-30048-7_10
 - [14] Rina Dechter, K. Kask, E. Bin, and R. Emek. 2002. Generating random solutions for constraint satisfaction problems. In *Proc. of AAAI*. 15–21.
 - [15] Joe W Duran and Simeon C Ntafos. 1984. An evaluation of random testing. *IEEE Transactions on Software Engineering* 4 (1984), 438–444. <https://doi.org/10.1109/TSE.1984.5010257>
 - [16] Rafael Dutra, Kevin Laeufer, Jonathan Bachrach, and Koushik Sen. 2018. Efficient sampling of SAT solutions for testing. In *2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE)*. IEEE, 549–559. <https://doi.org/10.1145/3180155.3180248>
 - [17] Stefano Ermon, Carla P Gomes, Ashish Sabharwal, and Bart Selman. 2013. Embed and Project: Discrete Sampling with Universal Hashing. In *Proc. of NIPS*. 2085–2093.
 - [18] Vibhav Gogate and Rina Dechter. 2006. A new algorithm for sampling CSP solutions uniformly at random. In *CP*. Springer, 711–715. https://doi.org/10.1007/11889205_56
 - [19] Loreto Gonzalez-Hernandez, Nelson Rangel-Valdez, and Jose Torres-Jimenez. 2010. Construction of mixed covering arrays of variable strength using a tabu search approach. In *International Conference on Combinatorial Optimization and Applications*. Springer, 51–64. https://doi.org/10.1007/978-3-642-17458-2_6
 - [20] Rahul Gupta, Shubham Sharma, Subhjit Roy, and Kuldeep S. Meel. 2019. WAPS: Weighted and Projected Sampling. In *Proceedings of Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*.
 - [21] Mahesh A. Iyer. 2003. RACE: A word-level ATPG-based constraints solver system for smart random simulation. In *Proc. of ITC*. 299–308. <https://doi.org/10.1109/TEST.2003.1270852>
 - [22] Mark R. Jerrum and Alistair Sinclair. 1996. The Markov chain Monte Carlo method: an approach to approximate counting and integration. *Approximation algorithms for NP-hard problems* (1996), 482–520.
 - [23] Mark R. Jerrum, Leslie G. Valiant, and Vijay V. Vazirani. 1986. Random generation of combinatorial structures from a uniform distribution. *Theoretical Computer Science* 43, 2-3 (1986), 169–188. [https://doi.org/10.1016/0304-3975\(86\)90174-X](https://doi.org/10.1016/0304-3975(86)90174-X)
 - [24] Martin Fagereng Johansen, Øystein Haugen, and Franck Fleurey. 2012. An algorithm for generating t-wise covering arrays from large feature models. In *Proceedings of the 16th International Software Product Line Conference-Volume 1*. 46–55. <https://doi.org/10.1145/2362536.2362547>
 - [25] Christian Kaltenecker, Alexander Grebhahn, Norbert Siegmund, Jianmei Guo, and Sven Apel. 2019. Distance-based sampling of software configuration spaces. In *Proceedings of the 41st International Conference on Software Engineering*. IEEE Press, 1084–1094. <https://doi.org/10.1109/ICSE.2019.00112>
 - [26] Richard M Karp. 1972. Reducibility among combinatorial problems. In *Complexity of computer computations*. Springer, 85–103. https://doi.org/10.1007/978-1-4684-2001-2_9
 - [27] Nathan Kitchen. 2010. *Markov Chain Monte Carlo Stimulus Generation for Constrained Random Simulation*. Ph.D. Dissertation. University of California, Berkeley.
 - [28] Nathan Kitchen and Andreas Kuehlmann. 2007. Stimulus generation for constrained random simulation. In *Proc. of ICCAD*. 258–265. <https://doi.org/10.1109/ICCAD.2007.4397275>
 - [29] Alexander Knüppel, Thomas Thüm, Stephan Mennicke, Jens Meinicke, and Ina Schaefer. 2017. Is there a mismatch between real-world feature models and product-line research?. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. ACM, 291–302. <https://doi.org/10.1145/3106237.3106252>
 - [30] D Richard Kuhn, Raghu N Kacker, and Yu Lei. 2010. Practical combinatorial testing. *NIST special Publication* 800, 142 (2010), 142.
 - [31] Yu Lei, Raghu Kacker, D Richard Kuhn, Vadim Okun, and James Lawrence. 2007. IPOG: A general strategy for t-way software testing. In *14th Annual IEEE International Conference and Workshops on the Engineering of Computer-Based Systems (ECBS'07)*. IEEE, 549–556. <https://doi.org/10.1109/ECBS.2007.47>
 - [32] Yu Lei, Raghu Kacker, D Richard Kuhn, Vadim Okun, and James Lawrence. 2008. IPOG/IPOG-D: efficient test generation for multi-way combinatorial testing. *Software Testing, Verification and Reliability* 18, 3 (2008), 125–148. <https://doi.org/10.1002/stvr.381>
 - [33] Yu Lei and Kuo-Chung Tai. 1998. In-parameter-order: A test generation strategy for pairwise testing. In *Proceedings Third IEEE International High-Assurance Systems Engineering Symposium (Cat. No. 98EX231)*. IEEE, 254–261.
 - [34] Jia Hui Liang, Vijay Ganesh, Krzysztof Czarnecki, and Venkatesh Raman. 2015. Sat-based analysis of large real-world feature models is easy. In *Proceedings of the 19th International Conference on Software Product Line*. ACM, 91–100. <https://doi.org/10.1145/2791060.2791070>
 - [35] Jinkun Lin, Chuan Luo, Shaowei Cai, Kaile Su, Dan Hao, and Lu Zhang. 2015. TCA: An efficient two-mode meta-heuristic algorithm for combinatorial test generation (t). In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 494–505. <https://doi.org/10.1109/ASE.2015.61>
 - [36] Roberto E Lopez-Herrejon, Francisco Chicano, Javier Ferrer, Alexander Egyed, and Enrique Alba. 2013. Multi-objective optimal test suite computation for software product line pairwise testing. In *2013 IEEE International Conference on Software Maintenance*. IEEE, 404–407. <https://doi.org/10.1109/ICSM.2013.58>
 - [37] Neal Madras. 2002. Lectures on Monte Carlo Methods, Fields Institute Monographs 16. *American Mathematical Society* (2002). <https://doi.org/10.1090/fim/016>
 - [38] Robert Mandl. 1985. Orthogonal Latin squares: an application of experiment design to compiler testing. *Commun. ACM* 28, 10 (1985), 1054–1058. <https://doi.org/10.1145/4372.4375>
 - [39] Dusica Marijan, Arnaud Gottlieb, Sagar Sen, and Aymeric Hervieu. 2013. Practical pairwise testing for software product lines. In *Proceedings of the 17th international software product line conference*. 227–235. <https://doi.org/10.1145/2491627.2491646>
 - [40] James D McCaffrey. 2009. Generation of pairwise test sets using a genetic algorithm. In *2009 33rd annual IEEE international computer software and applications conference*, Vol. 1. IEEE, 626–631. <https://doi.org/10.1109/COMPSAC.2009.91>
 - [41] Flávio Medeiros, Christian Kästner, Márcio Ribeiro, Rohit Gheyi, and Sven Apel. 2016. A comparison of 10 sampling algorithms for configurable systems. In *2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE)*. IEEE, 643–654. <https://doi.org/10.1145/2884781.2884793>
 - [42] Kuldeep S. Meel. 2014. *Sampling Techniques for Boolean Satisfiability*. Rice University. M.S. Thesis.
 - [43] George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. 1978. An analysis of approximations for maximizing submodular set functions. *Mathematical programming* 14, 1 (1978), 265–294. <https://doi.org/10.1007/BF01588971>
 - [44] Changhai Nie and Hareton Leung. 2011. A survey of combinatorial testing. *ACM Computing Surveys (CSUR)* 43, 2 (2011), 1–29. <https://doi.org/10.1145/1883612.1883618>
 - [45] Kari J Nurmela. 2004. Upper bounds for covering arrays by tabu search. *Discrete applied mathematics* 138, 1-2 (2004), 143–152. [https://doi.org/10.1016/S0166-218X\(03\)00291-9](https://doi.org/10.1016/S0166-218X(03)00291-9)
 - [46] Jeho Oh, Don Batory, Margaret Myers, and Norbert Siegmund. 2017. Finding near-optimal configurations in product lines by random sampling. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. 61–71. <https://doi.org/10.1145/3106237.3106273>
 - [47] Jeho Oh, Paul Gazzillo, Don Batory, Marijn Heule, and Maggie Myers. 2019. Uniform sampling from kconfig feature models. *The University of Texas at Austin, Department of Computer Science, Tech. Rep. TR-19-02* (2019).
 - [48] Jeho Oh, Paul Gazzillo, and Don S. Batory. 2019. t-wise coverage by uniform sampling. In *Proceedings of the 23rd International Systems and Software Product Line Conference, SPLC 2019, Volume A, Paris, France, September 9-13, 2019*. 15:1–15:4. <https://doi.org/10.1145/3336294.3342359>

- [49] Gilles Perrouin, Sagar Sen, Jacques Klein, Benoit Baudry, and Yves Le Traon. 2010. Automated and scalable t-wise test case generation strategies for software product lines. In *2010 Third international conference on software testing, verification and validation*. IEEE, 459–468. <https://doi.org/10.1109/ICST.2010.43>
- [50] Tobias Pett, Thomas Thüm, Tobias Runge, Sebastian Krieter, Malte Lochau, and Ina Schaefer. [n. d.]. Product Sampling for Product Lines: The Scalability Challenge. ([n. d.]).
- [51] Quentin Plazar, Mathieu Acher, Gilles Perrouin, Xavier Devroey, and Maxime Cordy. 2019. Uniform sampling of sat solutions for configurable systems: Are we there yet?. In *2019 12th IEEE Conference on Software Testing, Validation and Verification (ICST)*. IEEE, 240–251. <https://doi.org/10.1109/ICST.2019.00032>
- [52] Atri Sarkar, Jianmei Guo, Norbert Siegmund, Sven Apel, and Krzysztof Czarnecki. 2015. Cost-efficient sampling for performance prediction of configurable systems (t). In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 342–352. <https://doi.org/10.1109/ASE.2015.45>
- [53] Shubham Sharma, Rahul Gupta, Subhjit Roy, and Kuldeep S. Meel. 2018. Knowledge Compilation meets Uniform Sampling. In *Proceedings of International Conference on Logic for Programming Artificial Intelligence and Reasoning (LPAR)*.
- [54] Norbert Siegmund, Alexander Grebhahn, Sven Apel, and Christian Kästner. 2015. Performance-influence models for highly configurable systems. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. 284–294. <https://doi.org/10.1145/2786805.2786845>
- [55] Norbert Siegmund, Sergiy S Kolesnikov, Christian Kästner, Sven Apel, Don Batory, Marko Rosenmüller, and Gunter Saake. 2012. Predicting performance via automated feature-interaction detection. In *Proceedings of the 34th International Conference on Software Engineering*. IEEE Press, 167–177. <https://doi.org/10.1109/ICSE.2012.6227196>
- [56] Neil JA Sloane. 1993. Covering arrays and intersecting codes. *Journal of combinatorial designs* 1, 1 (1993), 51–63. <https://doi.org/10.1002/jcd.3180010106>
- [57] Reinhard Tartler, Daniel Lohmann, Christian Dietrich, Christoph Egger, and Julio Sincero. 2011. Configuration coverage in the analysis of large-scale system software. In *Proceedings of the 6th Workshop on Programming Languages and Operating Systems*. 1–5. <https://doi.org/10.1145/2039239.2039242>
- [58] Keizo Tatsumi. 1987. Test case design support system. In *Proc. International Conference on Quality Control (ICQC'87)*. 615–620.
- [59] Thomas Thüm, Sven Apel, Christian Kästner, Ina Schaefer, and Gunter Saake. 2014. A classification and survey of analysis strategies for software product lines. *ACM Computing Surveys (CSUR)* 47, 1 (2014), 1–45. <https://doi.org/10.1145/2580950>
- [60] Thomas Thum, Christian Kastner, Sebastian Erdweg, and Norbert Siegmund. 2011. Abstract features in feature modeling. In *2011 15th International Software Product Line Conference*. IEEE, 191–200. <https://doi.org/10.1109/SPLC.2011.53>
- [61] Ziyuan Wang, Baowen Xu, and Changhai Nie. 2008. Greedy heuristic algorithms to generate variable strength combinatorial test suite. In *2008 The Eighth International Conference on Quality Software*. IEEE, 155–160. <https://doi.org/10.1109/QSIC.2008.52>
- [62] Akihisa Yamada, Armin Biere, Cyrille Artho, Takashi Kitamura, and Eun-Hye Choi. 2016. Greedy combinatorial test case generation using unsatisfiable cores. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*. 614–624. <https://doi.org/10.1145/2970276.2970335>