



# Uniform and scalable sampling of highly configurable systems

Ruben Heradio<sup>1</sup> · David Fernandez-Amoros<sup>1</sup> · José A. Galindo<sup>2</sup> · David Benavides<sup>2</sup> · Don Batory<sup>3</sup>

Accepted: 9 December 2021

© The Author(s) 2022

## Abstract

Many analyses on configurable software systems are intractable when confronted with colossal and highly-constrained configuration spaces. These analyses could instead use statistical inference, where a tractable sample accurately predicts results for the entire space. To do so, the laws of statistical inference requires each member of the population to be equally likely to be included in the sample, i.e., the sampling process needs to be “uniform”. SAT-samplers have been developed to generate uniform random samples at a reasonable computational cost. However, there is a lack of experimental validation over colossal spaces to show whether the samplers indeed produce uniform samples or not. This paper (i) proposes a new sampler named **BDDSampler**, (ii) presents a new statistical test to verify sampler uniformity, and (iii) reports the evaluation of **BDDSampler** and five other state-of-the-art samplers: **KUS**, **QuickSampler**, **Smarch**, **Spur**, and **Unigen2**. Our experimental results show only **BDDSampler** satisfies both scalability and uniformity.

**Keywords** Uniform sampling · Configurable systems · Software product lines · Binary decision diagrams · SAT-solvers

---

Communicated by : Philippe Collet, Sarah Nadi and Christoph Seidl

This article belongs to the Topical Collection: *Software Product Lines and Variability-rich Systems (SPLC)*

---

✉ Ruben Heradio  
rheradio@issi.uned.es  
David Fernandez-Amoros  
david@issi.uned.es  
José A. Galindo  
jagalindo@us.es  
David Benavides  
benavides@us.es  
Don Batory  
batory@cs.utexas.edu

<sup>1</sup> Universidad Nacional de Educación a Distancia, Madrid, Spain

<sup>2</sup> University of Seville, Seville, Spain

<sup>3</sup> University of Texas at Austin, Austin, Texas, USA

# 1 Introduction

Generating random SAT-solutions is of critical importance in several domains: *Software Product Lines* (SPLs) analysis and configuration (Krieter 2019; Muñoz et al. 2019; Oh et al. 2017), software testing (Chakraborty and Meel 2019; Dutra et al. 2018; Plazar et al. 2019; Roy et al. 2018), and integrated circuit simulation and verification (Hübner and Becker 2011; Naveh et al. 2006; Yuan et al. 2002).

To get a sense of this problem's relevancy and complexity, consider an example taken from the SPL domain. **BusyBox**<sup>1</sup> is a software tool that replaces many standard GNU/Linux utilities with a single small executable, thus providing an environment customized for a diversity of embedded systems. To achieve size-optimization, **BusyBox** is remarkably modular, supporting the inclusion/exclusion of 613 features at compile time. These features and their interrelationships are specified with a configuration language named **Kconfig**.<sup>2</sup> To guarantee that every valid configuration satisfies all dependencies, the **Kconfig** model of **BusyBox** is translated into a Boolean formula that is then processed with a logic engine (Batory 2005; Fernandez-Amoros et al. 2019) (e.g., a SAT solver (Biere et al. 2009)). A valid configuration corresponds to a satisfiable assignment of the formula, also called, a *SAT solution* (Plazar et al. 2019) or a *witness* (Chakraborty and Meel 2019).

As a consequence of the inter-feature dependencies, the space of valid configurations ( $7.428 \cdot 10^{146}$ ) is a tiny portion of the whole configuration space ( $2^{613}$ ): only  $2.185 \cdot 10^{-36}\%$  of the possible configurations are valid (Heradio et al. 2020). Nevertheless, the population of valid configurations is still colossal. Those SPL analyses that examine every valid configuration are unscalable.

For instance, Halin et al. (2019) adopted an exhaustive strategy to test the **JHipster**<sup>3</sup> system, checking all its valid configurations. **JHipster** is a code generator for web applications with 45 selectable features that can produce a total of 26,256 valid configurations. Checking this modest configuration space with the **INRIA Grid'5000**<sup>4</sup> required 4,376 hours of CPU time ( $\sim 182$  days), and 5.2 terabytes of disk space.

Others have advocated approaching this and related problems via statistical inference (Alf  rez et al. 2019; Alves Pereira et al. 2020; Guo et al. 2018; Kaltenecker et al. 2019; Kolesnikov et al. 2019; Nair et al. 2017; Oh et al. 2017; Temple et al. 2016; Weckesser et al. 2018); that is, working with a tractable sample that predicts the results for the entire population. An essential requirement is that all samples be genuinely representative of the population (Kaplan 2012). In other words, each member of the population must be equally likely to be included in a sample. Authors often use the term *uniform random sampling* (Oh et al. 2017; Plazar et al. 2019; Sharma et al. 2018) for this idea.

A naive approach to get such a sample would (i) generate a random configuration set without considering feature dependencies, and then (ii) check with a logic engine if each configuration conforms to those dependencies. Unfortunately, and as mentioned above, feature dependencies shrink the configuration space extraordinarily, and so getting a single valid configuration randomly is extremely unlikely. As a result, more advanced algorithms generate valid and uniform random samples at a reasonable computational cost.

<sup>1</sup> <https://busybox.net/>

<sup>2</sup> <https://www.kernel.org/doc/Documentation/kbuild/kconfig-language.txt>

<sup>3</sup> <https://www.jhipster.tech/>

<sup>4</sup> <https://www.grid5000.fr/>

Verifying that these algorithms and their tools indeed generate genuine uniform samples is a challenge by itself, because it requires examining the consistency between sample *statistics* and their corresponding population *parameters* (e.g., how frequently a feature appears in a sample compared to its probability of being included in every valid configuration (Heradio et al. 2019)). As configuration spaces can be colossal, current procedures that certify a sampler's uniformity has the severe shortcoming of requiring gigantic sample sizes to estimate reliable statistics (Dutra et al. 2018; Achlioptas et al. 2018; Chakraborty and Meel 2019). Consequently, sampler uniformity has been checked only on miniature models so far, which is not convincing. Also, most uniformity procedures compute population parameters in a poorly scalable way (e.g., requiring calling a #SAT solver thousands of times (Plazar et al. 2019)).

This paper extends our paper in SPLC'20 (Heradio et al. 2020), where (i) a statistical test is formulated to reduce the sample size required for assessing a samplers' uniformity, and (ii) population parameters are computed with scalable algorithms we proposed in Heradio et al. (2019). The additional contributions of this present paper are:

1. A new sampler called **BDDSampler**, which is built upon a *Binary Decision Diagram* (BDD) (Bryant 1986) technology (see Section 3).
2. A new statistical test to validate a samplers' uniformity, reducing the sample size requirements even more than our previous test (see Section 4).
3. An experimental validation with our new test of **BDDSampler** and other five state-of-the-art samplers (**KUS** (Sharma et al. 2018), **QuickSampler** (Dutra et al. 2018), **Spur** (Achlioptas et al. 2018), **Smarch** (Oh et al. 2019), and **Unigen2** (Chakraborty et al. 2013; Chakraborty 2015)) on configuration models with up to 18,570 variables (see Section 5).
4. Experimental results show (i) our new statistical test needs the smallest sample size of all existing uniformity validation methods, and (ii) **BDDSampler** is the only sampler that satisfies both uniformity and scalability. Our software artifacts (**BDDSampler**, and the data and code scripts for replicating the experiments) are freely available at public repositories (see Section 8).

## 2 Related work

Before discussing related work, a terminological clarification is needed. In the machine learning, the term *sample* usually refers to a single data point (Chollet and Allaire 2018). However, in *inferential statistics*, a *sample* is typically a collection of *cases*, where the number of cases in the sample is the *sample size* (Chihara and Hesterberg 2011; Kaplan 2012). This paper adopts this latter terminology, and consequently, a *sample* is a set of configurations (i.e., a collection of SAT-solutions), whose cardinal is its sample size.

Here is additional standard statistical terminology that we will use in this paper. Inferential statistics aims to generalize the results obtained from a sample to the entire population. To do so, the most widespread approach, called *Null Hypothesis Significance Test* (NHST), quantifies the probability of obtaining the sample results conditioned on the assumption that a given *null hypothesis* ( $H_0$ ) is true (NHST fundamentals are explained in Chapter 13 of Kaplan (2012) and Chapter 3 of Vasishth and Broe (2011)). If such probability (named *p-value*) is less or equal than an established threshold (called the *significance level* ( $\alpha$ )) then  $H_0$  is rejected, and thus its alternative hypothesis  $H_a$  accepted. Otherwise,  $H_0$  is kept.

As Table 1 shows, two mistakes under this framework can be made due to unusual random samples: rejecting a true  $H_0$  (named *Type 1 error*), and failing to reject a false  $H_0$  (called *Type 2 error*). The expression  $1 - \beta$  is known as the test's *power*. The experimenter can adjust the Type 1 and 2 error probabilities through the thresholds  $\alpha$  and  $\beta$  (see Chapter 4 of Vasishth and Broe (2011)).

## 2.1 Uniform random samplers

The following sections summarize some of the most common strategies to generate uniform random samples for a model of a configuration space that is encoded as a Boolean formula  $\varphi$ .

### 2.1.1 Atomic mutations (QuickSampler)

QuickSampler<sup>5</sup> (Dutra et al. 2018) uses a heuristic to gain scalability by minimizing the number of calls to a constraint solver. It generates a random configuration without taking into account the formula constraints. This configuration often violates constraints and thus is unsatisfiable. So, QuickSampler calls the Z3 solver (de Moura and Bjørner 2008) to fix the configuration by finding a MAX-SAT-solution. Then, QuickSampler flips the value of each variable and calls again Z3 to get another valid configuration. The differences between the variable values of the original and flipped SAT configurations are called *atomic mutations*. By combining mutations, QuickSampler quickly generates new configurations without calling the solver as those configurations are usually legal (Dutra et al. 2018).

### 2.1.2 Hashing-based sampling (Unigen2)

Several techniques divide the space of SAT-solutions into small “cells” of approximately the same size using  $r$  independent hash functions. Accordingly, sampling is done by choosing a cell at random, and then getting a satisfying assignment for that cell using a SAT solver. A critical point of these techniques is determining the “right”  $r$  value. For instance, Bellare et al. (2000) showed that an  $r$  equal to the number of formula variables guarantees uniformity. However, Chakraborty et al. (2013) reported that such  $r$  does not scale in practice; in contrast,  $r = 3$  scales better and ensures near-uniformity. Unigen2<sup>6</sup> (Chakraborty 2015) develops these ideas further, giving stronger uniformity guarantees.

### 2.1.3 Counting-based sampling (KUS, Smarch, and Spur)

In Section 7.1.4 of Knuth (2009), Knuth showed how to accomplish uniform random sampling by subsequently partitioning the SAT-solution space on variable assignments, and then counting the number of solutions of the resulting parts. Again,  $\varphi$  be a Boolean formula of  $v$  variables  $x_1, x_2, \dots, x_v$ ; let  $\#SAT(\varphi)$  denote the number of solutions to  $\varphi$ ; and let  $r \in [0, 1]$  be a random number in the unit interval. Conceptually, the procedure works as follows: The number of solutions where  $x_1$  is true is counted, namely  $\#SAT(\varphi \wedge x_1)$ .  $x_1$

<sup>5</sup> [https://github.com/RafaelTupynamba/quick\\_sampler](https://github.com/RafaelTupynamba/quick_sampler)

<sup>6</sup> <https://bitbucket.org/kuldeepmeel/unigen>

**Table 1** Type 1 and 2 errors under the NHST framework

	$H_0$ is true in reality ( $H_0$ )	$H_0$ is false in reality ( $\neg H_0$ )
The decision inferred from the sample is “reject $H_0$ ” ( $R$ )	$\Pr(R H_0) = \alpha$ <i>Type 1 error</i>	$\Pr(R \neg H_0) = 1 - \beta$ <i>Power</i>
The decision inferred from the sample is “do not reject $H_0$ ” ( $\neg R$ )	$\Pr(\neg R H_0) = 1 - \alpha$	$\Pr(\neg R \neg H_0) = \beta$ <i>Type 2 error</i>

follows a Bernoulli distribution with probability  $p_1 = \frac{\#SAT(\varphi \wedge x_1)}{\#SAT(\varphi)}$ .  $x_1$  is assigned false if  $r \leq p_1$ , true otherwise. Suppose  $x_1$  is assigned false. Then,  $x_2$  follows a Bernoulli distribution with probability  $p_2 = \frac{\#SAT(\varphi \wedge \bar{x}_1 \wedge x_2)}{\#SAT(\varphi \wedge \bar{x}_1)}$ , and it would be randomly assigned. The procedure advances until the last variable  $x_v$  is assigned, and thus the random solution is completed.

The original algorithm by Knuth is specified on BDDs, as the probabilities required for all the possible SAT-solutions are computed just once with a single BDD traversal, and then reused every time a random configuration is generated. Oh (2017) reinvented Knuth’s algorithm and was the first to implement and apply it to SPL analyses. Since then, Knuth’s algorithm has been adapted to other *knowledge compilation* and *Davis-Putnam-Logemann-Loveland* (DPLL) (Davis et al. 1962) approaches. In particular, (i) the KUS<sup>7</sup> sampler (Sharma et al. 2018) substitutes BDDs with *deterministic-Decomposable Negation Normal Forms* (d-DNNFs), and (ii) Spur<sup>8</sup> (Achlioptas et al. 2018) and Smarch<sup>9</sup> (Oh et al. 2019) count SAT solutions with a #SAT-solver named sharpSAT (Thurley 2006).

### 2.1.4 New: BDDSampler, a scalable and uniform sampler

Section 3 describes a new sampler called BDDSampler, which is based on Knuth’s algorithm and implemented on top of the CUDD<sup>10</sup> library for BDDs.

According to the experimental results reported in Section 5, the only sampler that satisfies both scalability and uniformity is BDDSampler. More specifically, evidence shows that:

- BDDSampler, KUS, QuickSampler, and Spur are considerably faster than Smarch and Unigen2.
- In terms of uniformity, there are three types of samplers: (i) those that mostly fail to produce uniform samples (QuickSampler), (ii) those that usually work but from time to time generate non-uniform samples (KUS and Spur), and (iii) those that always produce uniform samples (BDDSampler, Smarch, and Unigen2).

<sup>7</sup> <https://github.com/meelgroup/KUS>

<sup>8</sup> <https://github.com/ZaydH/spur>

<sup>9</sup> [https://github.com/jeho-oh/Kclause\\_Smarch](https://github.com/jeho-oh/Kclause_Smarch)

<sup>10</sup> <https://github.com/vscosta/cudd>

## 2.2 Prior work on testing sampler uniformity

The following sections summarize the methods that have been devised to test the uniformity of a random sampler  $\mathbb{S}$ .

### 2.2.1 Method 1: Generate a massive sample with $\mathbb{S}$ , and compare it with another one obtained simulating an ideal uniform sampler

This is the most common technique in the literature (Achlioptas et al. 2018; Chakraborty 2015; Dutra et al. 2018; Plazar et al. 2019; Sharma et al. 2018). First, the total number  $n$  of SAT-solutions is counted for the Boolean formula  $\varphi$ , typically using a #SAT-solver. Having  $n$ , the generation of a uniform sample with size  $s$  is simulated as follows: imagine that numbers  $1, 2, \dots, n$  are put into a box; then,  $s$  numbers are sampled with replacement from the box, guaranteeing that the probability each number has to be extracted is  $\frac{1}{n}$ .

For example, JHipster encompasses 26,256 valid configurations (Halin<sup>n</sup> et al. 2019). Figure 1 shows the histogram of a sample ten times greater than the number of configurations ( $s = 26,256 \cdot 10$ ), which has been obtained sampling with replacement from the set  $\{1, 2, \dots, 26256\}$ . The x-axis depicts numbers' occurrences, i.e., there are numbers that appear 0, 1, ..., 27 times in the sample; the y-axis shows how frequent are those occurrences in the sample. As expected, most numbers appear ten times (see the red vertical line in Fig. 1), however, and due to randomness, some numbers appear more frequently than others.

Another sample with size  $s$  (whose value is quantified shortly), is then generated with sampler  $\mathbb{S}$ . For this sample, a counterpart histogram to Fig. 1 is obtained, representing how often solutions appear in that sample.

Finally, the uniformity of  $\mathbb{S}$  is verified by measuring the distance between both histograms, using, for instance, the Kullback-Leibler divergence (Achlioptas et al. 2018).

Unfortunately, this method has a severe limitation: it does not scale except for formulas with a small number of SAT-solutions because, to produce reliable results,  $s$  needs to be much larger than  $n$  (see Achlioptas et al. (2018); Dutra et al. (2018) for an explanation). For example, Dutra et al. (2018) propose  $s \geq 5n$ . As the number of solutions grows exponentially with the number of variables of  $\varphi$ , the method only works for the simplest models with just a few features.

### 2.2.2 Method 2: Assume the existence of a uniform sampler $\mathbb{U}$ , and compare the samples generated by both $\mathbb{S}$ and $\mathbb{U}$

Chakraborty and Meel (2019) proposed this method and implementation called *barbarik*.<sup>11</sup> The method makes a strong assumption: there is a sampler  $\mathbb{U}$  that is known to be uniform. Thus, two samples of the same size  $s$  are generated with  $\mathbb{S}$  and  $\mathbb{U}$  and, depending on the *distance* between the samples, i.e., on how similar they are, *barbarik* decides if  $\mathbb{S}$  is approximately uniform.

The key of the method is how to define “approximately” for reaching a balance between uniformity and sample size, i.e., for avoiding the large  $s$  that Method 1 requires. Two parameters, called tolerance  $\epsilon$  and intolerance  $\eta$ , adjust the definition of

<sup>11</sup> <https://github.com/meelgroup/barbarik>

“uniformity” to avoid the above problems. A sampler is uniform whenever the probability  $p_1, p_2, \dots, p_n$  of all  $n$  solutions is exactly  $\frac{1}{n}$ .

Barbarik relaxes this definition, proposing that a sampler is *additive almost-uniform* if  $p_1, p_2, \dots, p_n \in [\frac{1-\epsilon}{n}, \frac{1+\epsilon}{n}]$ . Moreover, a sampler is  $\eta$ -far from uniformity if

$$\left| \sum_{i=1}^n p_i - \frac{1}{n} \right| \geq \eta$$

Chakraborty and Meel claim that  $s$  depends on  $\epsilon$  and  $\eta$  exclusively, but not on  $n$ . In particular, they state that a uniformity test with *significance level*  $\alpha = 0.1$  (i.e., 0.9 probability of accepting the uniformity of a sampler when it is genuinely uniform) and *Type 2 error*  $\beta = 0.1$  (i.e., 0.9 probability of rejecting the uniformity of a sampler that is not uniform) is accomplished when  $\epsilon = 0.6$  and  $\eta = 0.9$ , requiring a sample size of 1, 729, 750. Unfortunately, they do not provide a detailed formal proof for these settings in Chakraborty and Meel (2019).

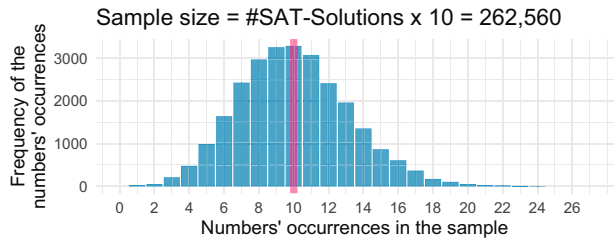
An evident weakness of this method is the necessity of a sampler  $\mathbb{U}$  with certified uniformity as a support lever. It is worth noting that, although an algorithm can be proven to generate uniform samples theoretically, some of its implementations may have errors. In other words, every sampling program needs to be tested, and thus Method 2 implicitly assumes the existence of another reliable uniformity testing method.

### 2.2.3 Method 3: Compare the theoretical variable probabilities in $\Pr(R|H_0) = \alpha$ with the empirical variable frequencies in a sample generated with $\mathbb{S}$

Plazar et al.’s method (Plazar et al. 2019) begins computing the theoretical probability each variable  $x$  has to appear in a SAT-solution. To do so, the procedure introduced in Section 2.1.3 is adopted, calling a #SAT solver repeatedly, one time per variable. #SAT( $\varphi$ ) gives the total number of SAT-solutions, and #SAT( $\varphi \wedge x$ ) calculates the number of solutions where  $x$  is true. Hence, the probability of  $x$  is  $p = \frac{\text{\#SAT}(\varphi \wedge x)}{\text{\#SAT}(\varphi)}$ . Likewise, if  $x$  is true  $t$  times in a sample of size  $s$ , its empirical frequency is  $f = \frac{t}{s}$ . Then, the *deviation* between  $p$  and  $f$  is  $d = 100 \cdot \frac{|p-f|}{p}$ . Finally, Plazar et al. propose two *thresholds* for  $d$ : (i) when  $d \leq 10$  for all variables, the deviations are *very low*, and thus sampler uniformity is accepted; (ii) when  $d \geq 50$  for some variables, they show *very high* deviations, and so uniformity is rejected. Regarding the sample size, Plazar et al. propose always using  $s \sim 10^6$ , independently of the number of variables of  $\varphi$  (no formal justification is given for this specific value in Plazar et al. (2019)).

Regrettably, this method often throws false negatives for variables with low probabilities. Suppose a variable has  $p = 0.01$ . Then, a genuine uniform sampler might easily generate a sample where  $f$  is slightly different just due to randomness, e.g.,  $f = 0.015$ . Therefore,  $d = 100 \cdot \frac{|0.01-0.015|}{0.01} = 50$ , and thus the sampler uniformity would be rejected. The chances that these types of wrong diagnoses happen increases with the number of low-probability variables, and it is worth noting that real models with numerous low-probability variables are not “corner cases”; for example, in three out of the seven configuration models analyzed in Heradio et al. (2019), more than 46% of their variables have  $p \leq 0.05$ : the open-source project *Fiasco v2014092821*, the *Dell* laptop configurator, and the *Automotive 02* system.

**Fig. 1** Simulated uniform random sample of the JHipster configuration model



#### 2.2.4 Method 4: a statistical goodness-of-fit test that compares theoretical variable probabilities in $\Pr(R|\neg H_0) = 1 - \beta$ with the empirical variable frequencies in a sample generated with $\mathbb{S}$

In the past (Heradio et al. 2020), we presented a procedure called *Feature Probability* (FP) test, which compares the empirical feature frequencies in a sample with the theoretical feature probabilities in the whole population of SAT-solutions. Instead of using the limited Method 3 *deviation* measure, our FP method (i) has a robust mathematical basis, (ii) estimates the *statistical significance* of the results (i.e., how generalizable they are), and (iii) supports adjusting the sample size according to precise statistical criteria (i.e., *Type 1* and *2 errors*, and *effect size*).

It is worth noting that a major shortcoming of Methods 1, 2, and 3 is the large sample size they need. For instance, in Achlioptas et al. (2018) and Sharma et al. (2018), Method 1 is applied on a model called `blasted_case110` with 287 variables, requiring  $s = 4 \cdot 10^6$  SAT-solutions. In Chakraborty and Meel (2019), Method 2 is used on `blasted_case110` as well, needing this time 1, 729, 750 SAT-solutions to ensure probability errors of Type 1  $\alpha = 0.1$  and Type 2  $\beta = 0.1$ . In contrast, our FP test provides stronger test guarantees ( $\alpha = 0.01$  and  $\beta = 0.01$ ) for `blasted_case110` with a minimal sample size of 13,027 solutions (i.e., a 99.25% sample size reduction with respect to Method 2).

#### 2.2.5 New: an improved goodness-of-fit test

Section 4 presents a new procedure that improves Method 4 by, instead of examining the variable probabilities, analyzing how the number of variables assigned to true distributes along the SAT-solutions. We show in Section 5 that the new method requires even smaller samples, thus widening the support for testing samplers' uniformity on larger models. For example, the sample size our new method requires for `blasted_case110` with  $\alpha = 0.01$  and  $\beta = 0.01$  becomes 6,563 solutions.

#### 2.2.6 Recap

Excluding the methods presented in this article and in our conference paper, there are serious practical problems in applying existing ways to test for sampler uniformity. We provide experimental evidence in Section 5 that our improved goodness-of-fit test is superior to prior work as it requires the smallest sample size of all existing tests, thus enabling the verification of samplers' uniformity in large models. As we will see, this highly increases the test sensitivity to detect samplers' uniformity flaws. Moreover, results show that our test



provides (i) *valid* judgements, which are consistent with the verdicts given by the alternative methods proposed in the literature, and (ii) *reliable* judgements, which remain consistent when the test is applied repeatedly to the same model and sampler.

### 3 The BDDSampler Tool

This section describes BDDSampler: a sampler that uses *Binary Decision Diagrams* (BDDs). A practical example how configuration models can be translated into Boolean formulas is presented in Section 3.1. Then, a BDD encoding of a Boolean formula is covered in Section 3.2. Finally, how BDDSampler works is explained in Section 3.3.

#### 3.1 From configuration models to Boolean formulas

Let us start with an example to help to explain BDDSampler and our samplers' uniformity test. As already mentioned in this paper's introduction, BusyBox supports the inclusion/exclusion of a number of features at compile time. These features and their interrelationships are specified with a configuration language named Kconfig which is used in many other relevant open-source projects (Berger et al. 2013), such as the Linux Kernel, axTLS, EmbToolkit, Freetz, etc.

Figure 2 shows an excerpt of the Kconfig specification of BusyBox v1.23.2. There are several configs encoding six features and their interdependencies. All features (STATIC, PIE, ..., FEATURE\_SHARED\_BUSYBOX) are Boolean (see the bool keyword in Lines 2, 4, ..., 15), meaning that they can be either selected or deselected. Configs trigger a prompt to request the user for their Boolean feature value, e.g., Build BusyBox as a static binary (no shared libs) in Line 2. Finally, some dependencies between features are set, e.g., according to the depends sentence in Line 10, BUILD\_LIBBUSYBOX can only be selected if none of the following features are selected: FEATURE\_PREFER\_APPLETS, PIE, neither STATIC.

The graph in Fig. 3 depicts the entire BusyBox configuration model, which includes 613 features and 530 inter-dependencies; nodes represent features, and edges depict dependencies. The Kconfig excerpt in Fig. 2 is zoomed in Fig. 3.

Given the configuration models' complexity, they are usually translated into Boolean formulas that are then processed with logic engines. For instance, Eq. 1 is the Boolean encoding of Fig. 2 (a detailed explanation of how to convert Kconfig specifications into Boolean formulas is given in Fernandez-Amoros et al. (2019)). In this section and the following one, we explain how to use BDDs for (i) generating random samples from the formulas, and (ii) testing the uniformity of an input sampler.

$$\begin{aligned} \varphi \equiv & (\neg \text{STATIC} \vee \neg \text{PIE}) \wedge \\ & (\neg \text{BUILD\_LIBBUSYBOX} \vee \neg \text{FEATURE\_PREFER\_APPLETS}) \wedge \\ & (\neg \text{BUILD\_LIBBUSYBOX} \vee \neg \text{PIE}) \wedge \\ & (\neg \text{BUILD\_LIBBUSYBOX} \vee \neg \text{STATIC}) \wedge \\ & (\neg \text{FEATURE\_INDIVIDUAL} \vee \text{BUILD\_LIBBUSYBOX}) \wedge \\ & (\neg \text{FEATURE\_SHARED\_BUSYBOX} \vee \text{BUILD\_LIBBUSYBOX}) \end{aligned} \quad (1)$$

### 3.2 A brief introduction to BDDs

A BDD (Bryant 1986) encodes a Boolean formula as a rooted directed acyclic graph composed of terminal and non-terminal nodes. Terminal nodes are represented as  $\boxed{0}$  and  $\boxed{1}$ , and non-terminal nodes are labeled with the formula variables. Two edges, named *low* and *high*, come out of every non-terminal node. Low is depicted with a dashed line ( $\dashrightarrow$ ), and high with a solid line ( $\rightarrow$ ). A BDD encodes every possible assignment of the formula variables as a path that descends from the root to the terminal nodes, going through solid lines when the corresponding variables are assigned to true and through dashed lines otherwise. An assignment is satisfiable, i.e., it evaluates the formula to true, whenever the traversed path ends at  $\boxed{1}$ .

Figure 4 depicts a BDD that encodes the **BusyBox** excerpt specified by Eq. 1. A configuration whose only activated features are `BUILD_LIBBUSYBOX` and `FEATURE_SHARED_BUSYBOX` conforms with the constraints (i.e., it is valid) and so it corresponds to the path `BUILD_LIBBUSYBOX`  $\rightarrow$  `FEATURE_INDIVIDUAL`  $\dashrightarrow$  `FEATURE_SHARED_BUSYBOX`  $\rightarrow$  `FEATURE_PREFER_APPLETS`  $\dashrightarrow$  `STATIC`  $\dashrightarrow$  `PIE`  $\dashrightarrow$   $\boxed{1}$ . In contrast, as `STATIC` and `PIE` are mutually exclusive, no configuration includes them simultaneously. Thus, all paths with solid lines coming out of both `STATIC` and `PIE` finish at  $\boxed{0}$ .

BDDs are typically ordered and reduced. A BDD is *ordered* when its variables are in the same position, called *index*, in every path from the root to the terminal nodes. For example, in Fig. 4, `STATIC` (whose index is 4) always goes before `PIE` and after `FEATURE_PREFER_APPLETS` (whose indices are 5 and 3, respectively). A BDD is *reduced* if it is free of redundant information. For instance, every blue/dark-shaded node in Fig. 4 is superfluous because both of its edges point to the same node and thus the formula evaluation is identical whether these variables are assigned to true or false. Consequently, these unnecessary tests are avoided in the reduced BDD in Fig. 5 to save computer memory.

It is worth noting that the variable ordering chosen to build the BDD has a tremendous impact on its size. Whereas a BDD can be reduced optimally (the reduction procedure was presented in the seminal article (Bryant 1986)), obtaining the best variable arrangement that minimizes its size is an NP-problem (Chapters 8 and 9 of Meinel and Theobald (1998) provide a comprehensive discussion on this topic). Several variable ordering heuristics (Fernandez-Amoros et al. 2020; Fernandez-Amoros et al. 2019; Mendonça 2009; Narodytska and Walsh 2007) have been proposed for the specific case of configuration model formulas. As reported in Section 5, we have been able to synthesize BDDs for large configuration models, with up to 17,000 features, by using these heuristics.

### 3.3 How BDDSampler works

BDDSampler takes an ordered and reduced BDD as input and generates random configurations in a two-step process described by Algorithms 1 and 2. Figure 5 summarizes Algorithm 1 computations for our running example. The algorithm decorates each non-terminal node with its probability of reaching the terminal  $\boxed{1}$  if the associated variable is set to true. Algorithm 1 proceeds in a bottom-up fashion, collecting the number of SAT solutions that can be produced by its low and high children (`solLO` and `solHI` in Lines 8-9), adding them up (`sol` in Line 10), and then computing the ratio corresponding to the high child (`pr` in Line 11). As the BDD is reduced, Algorithm 1 adjusts the solution counts in Lines 8-9 for the removed nodes with the expression  $2^{\text{index}(n_{\text{LOHI}}) - \text{index}(n) - 1}$ . For traversing efficiently the BDD,

```

1 config STATIC
2   bool "Build BusyBox as a static binary (no shared libs)"
3 config PIE
4   bool "Build BusyBox as a position independent executable"
5   depends on !STATIC
6 config FEATURE_PREFER_APPLETS
7   bool "exec prefers applets"
8 config BUILD_LIBBUSYBOX
9   bool "Build shared libbusybox"
10  depends on !FEATURE_PREFER_APPLETS && !PIE && !STATIC
11 config FEATURE_INDIVIDUAL
12   bool "Produce a binary for each applet, linked against libbusybox"
13  depends on BUILD_LIBBUSYBOX
14 config FEATURE_SHARED_BUSYBOX
15   bool "Produce additional busybox binary linked against libbusybox"
16  depends on BUILD_LIBBUSYBOX

```

Fig. 2 Excerpt of the BusyBox Kconfig specification

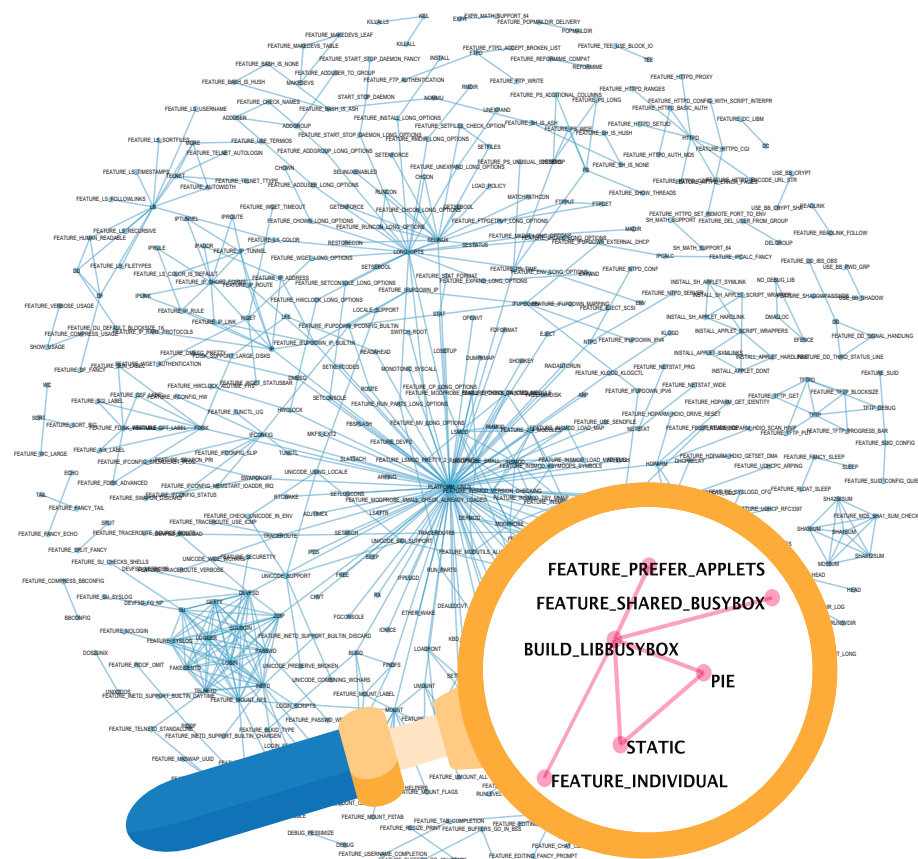
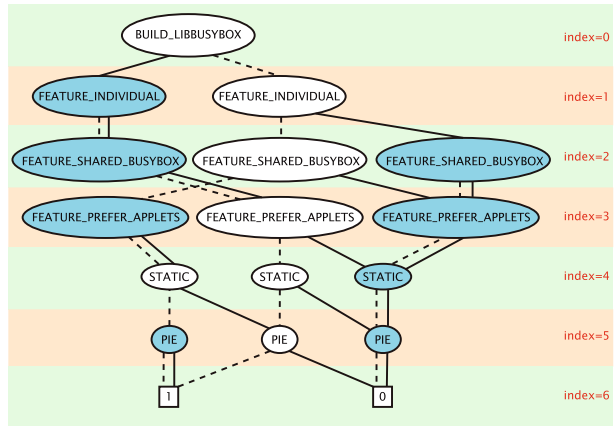
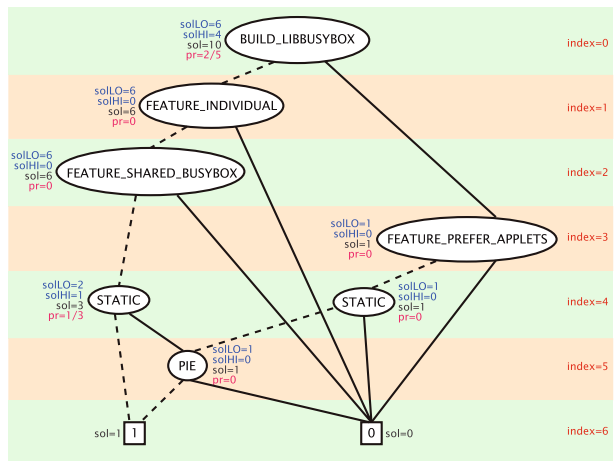


Fig. 3 Graph-representation of the BusyBox Kconfig specification

**Fig. 4** *Non-reduced* BDD encoding of Eq. 1



**Fig. 5** *Reduced* BDD encoding of Eq. 1



Algorithm 1 uses Bryant's method (Bryant 1986) as follows: the algorithm is called in Line 12 with the BDD root as argument and with a Boolean *mark* for every node being either all true or all false; then, it explores all nodes by recursively visiting the low and high children (Lines 6 and 7). Whenever a node is visited, its mark value is complemented (Line 2). Comparing the node with its children's marks, it is decided if the children have already been visited. The method ensures that each node is visited exactly once and that, when the traverse finishes, all node marks have the same value.

**Algorithm 1** Get all node probabilities.

---

```

1 Function getNodePr (n)
2   mark(n)  $\leftarrow \neg$  mark(n)
3   if n = 0 then sol[n]  $\leftarrow$  0 // [0] is reached
4   else if n = 1 then sol[n]  $\leftarrow$  1 // [1] is reached
5   else
6     // explore low
7     if mark(n)  $\neq$  mark( $n_{LO}$ ) then getNodePr ( $n_{LO}$ )
8     // explore high
9     if mark(n)  $\neq$  mark( $n_{HI}$ ) then getNodePr ( $n_{HI}$ )
10    // get node probabilities
11    solLO  $\leftarrow$  sol[ $n_{LO}$ ]  $\cdot 2^{\text{index}(n_{LO}) - \text{index}(n) - 1}$ 
12    solHI  $\leftarrow$  sol[ $n_{HI}$ ]  $\cdot 2^{\text{index}(n_{HI}) - \text{index}(n) - 1}$ 
13    sol[n]  $\leftarrow$  solLO + solHI
14    pr[n]  $\leftarrow \frac{\text{solHI}}{\text{sol[n]}}$ 
15  getNodePr (ROOT)

```

---

Whereas Algorithm 1 is run once as an initialization method, Algorithm 2 needs to be run as many times as configurations we want to generate. Algorithm 2 performs a random walk from the root to the terminal [1]. When a non-reduced node is visited, the path is selected randomly according to its probability (Lines 11-16): if the node probability is  $p$ , then its low and high edges are chosen with probabilities  $1 - p$  and  $p$ , respectively. Regarding the reduced nodes, the generated configuration will be valid no matter if their variables are set to true or false (that is the reason why these nodes were removed). Thus their value is chosen randomly with a 1/2

**Algorithm 2** Generate a random Configuration.

---

```

1 Function random
2   return a random number  $\in [0, 1)$ 
3 Function fiftyfifty
4   return random() < 0.5
5 Function generateConfiguration()
6   // generate random values for variables corresponding to reduced
7   // ROOT predecessor nodes
8   for (i  $\leftarrow$  0; i < index(ROOT); i++) do
9     sample[i]  $\leftarrow$  fiftyfifty()
10  // generate random values for the remaining variables
11  trav  $\leftarrow$  ROOT
12  while trav  $\neq$  1 do // iterate until reaching the 1-terminal node
13    ind  $\leftarrow$  index(trav)
14    if random() < pr[trav] then
15      trav  $\leftarrow$  travHI
16      sample[ind]  $\leftarrow$  TRUE
17    else
18      trav  $\leftarrow$  travLO
19      sample[ind]  $\leftarrow$  FALSE
20  // generate random values for variables of reduced intermediate
21  // nodes
22  for (i  $\leftarrow$  index+1; i < index(trav); i++) do
23    sample[i]  $\leftarrow$  fiftyfifty()
24  return sample

```

---

probability by taking into account that a reduced node index may be less than the BDD root index (Lines 6-7) or greater (Lines 17-18).

Algorithm 2 is remarkably fast since its time complexity is proportional to the number of indices (i.e., of variables), not the number of nodes in the BDD. Moreover, multiple instances of Algorithm 2 can be run in parallel over the same BDD, as the node probabilities are read but not modified.

Finally, **BDDSampler** is built on top of CUDD 3.0.<sup>12</sup> As other modern BDD libraries like **Sylvan** (van Dijk 2016), CUDD uses a technique called *complement edges* (Brace et al. 1990) to save nodes. With this technique, edges are enriched with a complement attribute that removes the need of having two terminal-nodes (basically, when an edge has the complement attribute enabled, the only terminal node is interpreted as its negation). Accordingly, **BDDSampler** tweaks Algorithms 1 and 2 to work with complement edges. We have decided to show the algorithms for regular BDDs without complement arcs for simplicity.

The BDDs of 218 models, which will be used in Section 5 to perform our experimental evaluation, are available at <https://doi.org/10.5281/zenodo.4514919> in the DDDMP format that CUDD uses for complement edge BDDs.

## 4 Assessing the uniformity of SAT solution samplers

Figure 6 sketches our approach to verify that a sampler generates uniform random samples of a model that is encoded as a Boolean formula. The method compares empirical information about a sample with theoretical information about the whole population of SAT-solutions that the model represents.

### 4.1 The SFpC goodness-of-fit test

In statistics, the procedures for examining how well a sample agrees with the population distribution are known as *goodness-of-fit* tests (D'Agostino and Stephens 1986). They require characterizing both the sample and the population in terms of a quantitative measure. In particular, we propose the distribution of the number of variables assigned to true among all SAT-solutions, called the *Selected Features per Configuration* (SFpC) test. For instance, Fig. 7 compares the theoretical distribution of all  $7.428 \cdot 10^{146}$  SAT-solutions of the **BusyBox** model with the distribution of 17,738 configurations generated with the samplers **BDDSampler** and **QuickSampler**, Fig. 7a,b respectively. The justification for this sample size 17,738 is given in Section 4.2.

The distribution of the whole population of SAT-solutions of a model can be computed with the *Product Distribution* (PD)<sup>13</sup> algorithm we proposed in Heradio et al. (2019). PD takes the BDD encoding of a model as input, and as explained in Section 3.D of Heradio et al. (2019), its time complexity is  $O(nv^2)$ , where  $n$  is the number of BDD nodes and  $v$  the number of model variables. Accordingly, PD scales for large models. For instance, on an Intel(R) Core(TM) i7-6700HQ, it took 2.74 minutes to compute the distribution of the **Automotive02** model (Krieter et al. 2018), which with 17,365 variables and 321,897 clauses encompasses  $5.26 \cdot 10^{1,441}$  SAT-solutions.

<sup>12</sup> <https://github.com/vscosta/cudd>

<sup>13</sup> <https://github.com/rheradio/VMStatAnal>

As the theoretical histogram shows in Fig. 7a, the smallest and largest BusyBox configurations have 6 and 571 features activated, respectively. 95% of the configurations have between 277 and 327 variables assigned to true.

The BDDSampler histogram (Fig. 7a) agrees with the normally distributed population. However, the QuickSampler histogram (Fig. 7b) is bimodal where most configurations have 100 or 200 features approximately, quite different from the theoretical histogram.

After exploring the sample's goodness-of-fit graphically, it is desirable to advance towards a more formal test that provides an accurate numerical quantification. A good candidate to measure the distance/difference between the sample and population distributions is the *Kullback–Leibler divergence*<sup>14</sup> (Cover and Thomas 2006). For discrete probability distributions  $P$  and  $F$  specified on the same probability space  $\mathbb{X}$ , the Kullback–Leibler divergence from  $F$  to  $P$  is defined as:

$$D_{\text{KL}}(P||F) = \sum_{x \in \mathbb{X}} P(x) \log_2 \left( \frac{P(x)}{F(x)} \right) \quad (2)$$

However, the Kullback–Leibler divergence is not symmetric, and thus it cannot rigorously be considered a *metric* (Lin 1991). For this reason, we use its symmetrical and normalized version, which is named *Jensen-Shannon divergence* (Cover and Thomas 2006; Lin 1991) and defined as:

$$JSD(P||F) = \frac{1}{2} D_{\text{KL}}(P||M) + \frac{1}{2} D_{\text{KL}}(F||M) \quad (3)$$

where  $M = \frac{1}{2}(P + F)$ .

In our case, vectors  $F$  and  $P$  are defined as follows:

- $F = [f_0, f_1, \dots, f_n]$  stores the SAT-solution *frequency* distribution (i.e., the red histograms in Fig. 7). That is,

$$f_0 = \frac{\text{\#SAT solutions in the sample with no variables assigned to true}}{\text{sample size}}$$

$$f_1 = \frac{\text{\#SAT solutions in the sample with 1 variable assigned to true}}{\text{sample size}}$$

...

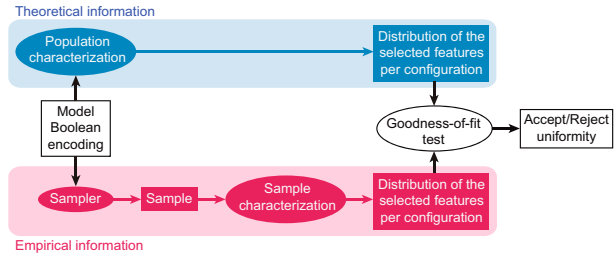
$$f_n = \frac{\text{\#SAT solutions in the sample with all variables assigned to true}}{\text{sample size}}$$

- $P = [p_0, p_1, \dots, p_n]$  stores the theoretical SAT-solution *probability* distribution of Fig. 7. That is,

$$p_i = \frac{\text{\#SAT solutions in the population with } i \text{ variables assigned to true}}{\text{population size}}$$

<sup>14</sup> The Kullback–Leibler divergence and especially one simplified version called *cross-entropy* are widely used as *loss functions* to compare the neural network predicted output with the observations used to train the network (see Chapter 3 of Goodfellow et al. (2016) for a summary of the Kullback–Leibler divergence and cross-entropy applications to deep learning).

**Fig. 6** Proposed method for verifying if a sampler generates uniform samples for a model



To avoid worthless comparisons, all  $i$ -elements with  $p_i = 0$  are removed from  $F$  and  $P$  because, as all solutions in the sample are guaranteed to be valid, the corresponding  $f_i$ 's are necessarily 0 as well. For instance, the **BusyBox** model has 613 variables, but all valid configurations have between 6 and 571 variables assigned to true. Therefore,  $\{f_0, f_1, \dots, f_5, f_{572}, f_{573}, \dots, f_{613}\}$  and  $\{p_0, p_1, \dots, p_5, p_{572}, p_{573}, \dots, p_{613}\}$  are deleted from  $F$  and  $P$ , respectively.

The Jensen-Shannon divergence  $\text{JSD}(P||F)$  measures to what extent the difference between  $F$  and  $P$  is greater than expected by chance if  $F$  corresponded to a uniform random sample. In the extreme cases,  $\text{JSD}(P||F) = 0$  when  $F$  totally matches  $P$ , and  $\text{JSD}(P||F) = 1$  when the  $F$  completely disagrees with  $P$ .

Nevertheless, JSD is a mere distance/difference metric, i.e., we cannot tell if JSD is *significantly* greater than expected due to randomness. Therefore, a *statistical inference test* is needed to quantify how generalizable the obtained distance is, i.e., a test that estimates the probability of a specific value of  $\text{JSD}(P||F)$  assuming that the sampler is genuinely uniform. In the case that the estimated probability is excessively low (below a significance level  $\alpha$ ), it is unlikely that the disagreement between  $F$  and  $P$  is due to chance, and so we can conclude that the sampler is not uniform.

Let  $s$  be the sample size (whose value we compute in Section 4.2), and  $m$  the number of elements in  $P$  after having removed those with  $p_i = 0$ . According to the proof given by Grosse et al. in Section 4.C of Grosse et al. (2002),  $2s(\ln 2)\text{JSD}(P||F)$  has a  $\chi^2$  distribution with  $m - 1$  degrees of freedom. As a result, a *Chi-Squared goodness-of-fit test* built upon the statistic  $2s(\ln 2)\text{JSD}(P||F)$  guides us to decide whether the sampler is uniform. In our **BusyBox** running example,  $s = 17,738$  and  $m = 613 - 6 - 42 = 565$ , hence if the sampler is uniform then  $2 \cdot 17,738(\ln 2)\text{JSD}(P||F)$  should follow a  $\chi_{565}^2$  distribution.

In contrast to typical *Null Hypothesis Significance Tests* (NHSTs), where the null hypothesis  $H_0$  states the opposite to what the researcher pursues to demonstrate, goodness-of-fit tests are a special case of NHSTs where  $H_0$  is: “the sample agrees with the population” (see Chapter 3 of D’Agostino and Stephens (1986) for a detailed description of *Chi-Squared goodness-of-fit tests*). Coming back to our case study, let us set the threshold  $\alpha = 0.01$  to test the **BusyBox** samples generated with:

– **BDDSampler**:

$$\text{JSD}(P||F) = 0.001085388$$

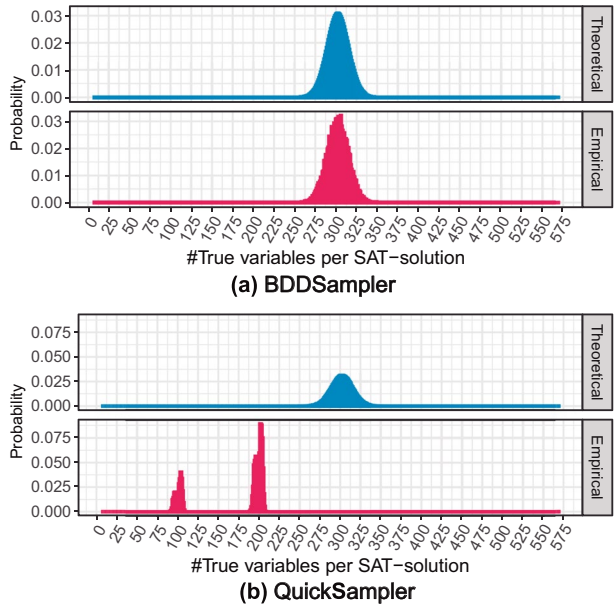
$$2 \cdot 17,738(\ln 2)\text{JSD}(P||F) = 26.68979$$

$$p\text{-value} = \Pr\left(\text{getting a value} \geq 26.68979 | H_0\right) \sim 1 > \alpha$$

⇒ Test result : Do not reject  $H_0$



**Fig. 7** Distribution of all BusyBox SAT-solutions compared with the distribution of 17,738 configurations generated with BDDSampler and QuickSampler



– QuickSampler:

$$\text{JSD}(P||F) \sim 1$$

$$2 \cdot 17,738(\ln 2)\text{JSD}(P||F) = 12,923.04$$

$$p\text{-value} = \Pr(\text{getting a value} \geq 12,923.04 | H_0) \sim 0 \leq \alpha$$

$\Rightarrow$  Test result : Reject  $H_0$

To sum up, the test corroborates numerically the histogram comparison in Fig. 7: BDDSampler generated a uniform sample, but QuickSampler did not.

## 4.2 Sample size estimation

The reliability of a Chi-Squared goodness-of-fit test depends on the following parameters (see Table 1):

- The *significance level*  $\alpha$  sets the probability of making a *Type 1 error*, i.e., the probability of rejecting  $H_0$  when it is indeed true (false positive). It is worth noting that  $\alpha$  is also the threshold for rejecting  $H_0$  (i.e.,  $H_0$  is rejected whenever the  $p$ -value  $\leq \alpha$ ).
- $\beta$  sets the probability of making a *Type 2 error*, i.e., the probability of accepting a false  $H_0$  (false negative). The expression  $1 - \beta$  is called the test's *power*, i.e., the probability of rejecting a false  $H_0$ .

When  $H_0$  is false, it is *false to some degree*. That degree is measured by another parameter called the *effect size* (Lakens 2013). In particular, Cohen (1988) proposes the index  $w$  for measuring the effect size in Chi-Squared tests. As a rule of thumb,  $w$  values of 0.1, 0.3, and 0.5 correspond to *small*, *medium*, and *large* effect sizes, respectively.

Interestingly, sample size, effect size,  $\alpha$ , and  $\beta$  have an intimate relationship in NHSTs: given any three of them, the fourth can be determined. In Section 7.3 of Cohen (1988), Cohen provides different *power tables* to estimate the minimum sample size required to ensure the reliability of a Chi-Squared test given the values of  $\alpha$ ,  $\beta$ ,  $w$ , and  $\chi^2$ 's degrees of freedom. Nowadays, there is available statistical software that provides those tables, e.g., the R package `pwr`<sup>15</sup> (see Chapter 10 of Kabacoff (2011)) and the `G*Power`<sup>16</sup> tool (Faul et al. 2007).

In the previous section, we saw that the goodness-of-fit of any sample from the `Bus-yBox` configuration model can be undertaken with a Chi-Squared test with 565 degrees of freedom. Then, according to Cohen's power tables, the required sample size is 17,738 configurations when  $\alpha = \beta = 0.01$ , and  $w = 0.1$ .

## 5 Empirical evaluation

This section describes the experimental evaluation of our approach using the *Goal/ Question/Metric* (GQM) method (van Solingen and Berghout 1999). As Fig. 8 shows, the evaluation pursues two goals (G1 and G2), which are refined into five questions (Q1-Q5) that are answered using different metrics.

The following points summarize our evaluation's goals and questions:

**G1: Samplers' evaluation.** The first goal G1 is to evaluate the scalability and uniformity of `BDDSampler` and the following state-of-the-art samplers: `KUS`<sup>17</sup> (Sharma et al. 2018), `QuickSampler`<sup>18</sup> (Dutra et al. 2018), `Smarch`<sup>19</sup> (Oh et al. 2019), `Spur`<sup>20</sup> (Achlioptas et al. 2018), and `Unigen2`<sup>21</sup> (Chakraborty et al. 2013; Chakraborty 2015). G1 is broken down into Questions Q1 and Q2:

**Q1: Samplers scalability.** *Are BDDSampler, KUS, QuickSampler, Smarch Spur, or Unigen2 able to generate samples with 1,000 configurations for models of all sizes within one hour?*

**Q2: Samplers' uniformity.** *Do BDDSampler, KUS, QuickSampler, Smarch Spur, or Unigen2 always generate uniform samples?*

**G2: SFpC's evaluation** The second goal G2 is to evaluate the scalability and quality, in terms of validity and reliability, of our SFpC test. G2 is refined into Questions Q3-Q5:

**Q3: SFpC's scalability.** *How much time and how many configurations does SFpC need to check the uniformity of a sampler on a model?*

**Q4: SFpC's validity.** *Does SFpC produce results consistent with the results obtained by other uniformity testing methods?*

**Q5: SFpC's reliability.** *When SFpC is applied repeatedly to the same model and sampler, are the results consistent?*

<sup>15</sup> <https://cran.r-project.org/web/packages/pwr>

<sup>16</sup> <https://www.psychologie.hhu.de/arbeitsgruppen/allgemeine-psychologie-und-arbeitspsychologie/gpower.html>

<sup>17</sup> <https://github.com/meelgroup/KUS>

<sup>18</sup> <https://github.com/RafaelTupynamba/quicksampler>

<sup>19</sup> <https://github.com/jeho-oh/Smarch>

<sup>20</sup> <https://github.com/ZaydH/spur>

<sup>21</sup> <https://bitbucket.org/kuldeepmeel/unigen>

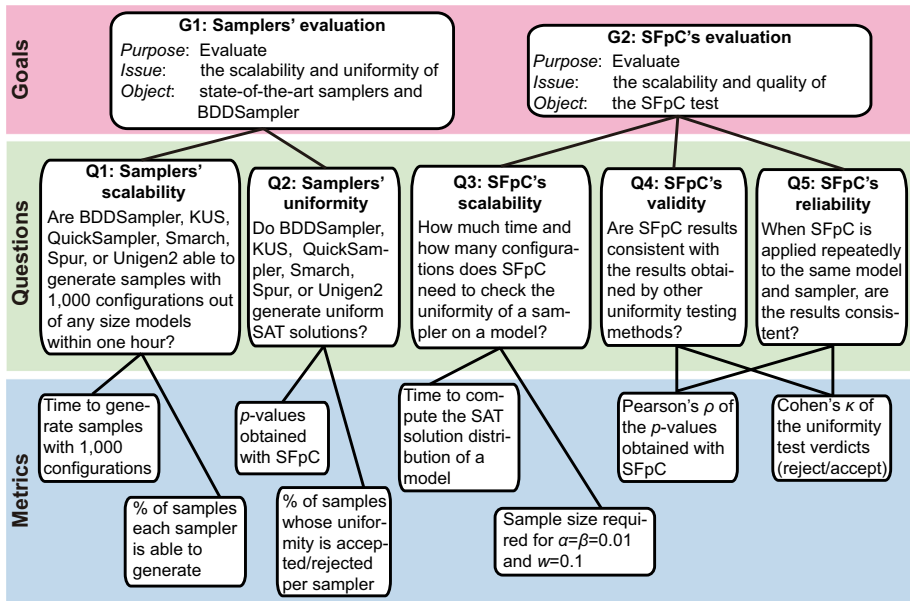


Fig. 8 Overview of the performed empirical evaluation with the GQM method

Section 5.1 presents the experimental setup. As Fig. 9 shows, three experiments E1-E3 were performed to solve the questions (e.g., Experiment E2 supported answering Questions Q2, Q3, and Q5). Sections 5.2-5.6 describe these experiments and the specific metrics used to answer the questions. The detailed results and all the material needed to replicate our experiments are available in the public repositories presented in Section 8.

## 5.1 Experimental setup

The samplers were tested against a suite of 218 models encoded as Boolean formulas in all of the following formats:

1. *DIMACS*, which is the format QuickSampler, Smarch, Spur, and Unigen2 use as input. These samplers rely on SAT technology, and DIMACS is the format for *Conjunctive Normal Form* (CNF) formulas that SAT technology uses.

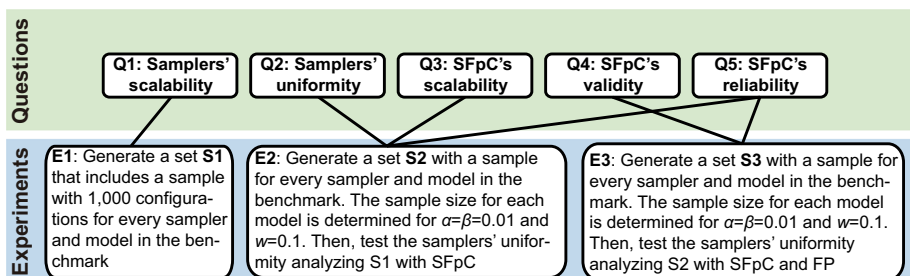


Fig. 9 Relationship between questions and experiments

2. *DDDMP*, which is the format **BDDSampler** and its underlying library **CUDD** use for BDDs.
3. *NNF*, which is the format **KUS** uses for d-DNNFs.

In particular,

- The DIMACS files of the industrial SAT formulas and **JHipster** were retrieved from Plazar et al. (2019).
- The DIMACS file of **LargeAutomotive** was gathered from Krieter et al. (2018).
- The DIMACS file of **DellSPLOT** was obtained from Nöhrer and Egyed (2013).
- We generated the DIMACS files of **axTLS**, **Fiasco**, **uClibc**, **ToyBox**, **BusyBox**, and **EmbToolkit** by processing their Kconfig specifications with our tool **Kconfig2Logic**<sup>22</sup> (Fernandez-Amoros et al. 2019).
- We generated all DDDMP files from their corresponding DIMACS files with our tool **Logic2BDD**<sup>23</sup> (Fernandez-Amoros et al. 2020).
- We generated all NNF files from their respective DIMACS files with the d-DNNF compiler **d4** (Lagniez and Marquis 2017) that is embedded in **KUS**.

In total, 209 are industrial SAT formulas (mostly modeling integrated circuits) that are typically used as a benchmark in the SAT-sampling literature (Achlioptas et al. 2018; Chakraborty 2015; Plazar et al. 2019). The remaining nine models represent configurable software systems. Table 2 describes the nine configuration models (the largest model is also referred as **Automotive02** in the SPL literature (Krieter et al. 2018)).

The histogram in Fig. 10 shows the model size distribution according to their number of variables. Since there is a wide range from the smallest model in the benchmark to the largest one (from 14 to 18,570 variables), the scale has been logarithmically transformed to shrink the range and thus facilitate the figure interpretation (see Chapter 5 of Winter (2020) for an explanation on logarithmic scale transformations). The scatter plot in Fig. 10 represents the model sizes in terms of their variables and clauses. The grey regression line shows that  $\text{Log}_2(\#\text{Clauses})$  depends on  $1.35 + 1.03 \cdot \text{Log}_2(\#\text{Variables})$ . Points corresponding to configuration models are labeled, and models with more and fewer clauses than those predicted by the linear regression are colored red and blue, respectively. Note that in the interval  $[9.88, 11.2]$  of  $\text{Log}_2(\#\text{Variables})$  there are only 5 models, and all of them have fewer clauses than predicted. As these models are simpler in terms of clauses, processing them requires less time than expected for their variable number, and thus regression curves in Figs. 12, 17, and 16 will show positive convexity in that interval.

The experiments were run on an Intel(R) Core(TM) i7-6700HQ, 2.60GHz, 16GB RAM, operating Linux Ubuntu 19.10. Samplers were executed on a single thread (i.e., with no parallelization), and without considering any Boolean formula preprocessing, such as *Minimal Independent Support* (MIS) (Ivrii et al. 2016).

<sup>22</sup> <https://github.com/davidfa71/Extending-Logic/tree/master/code/Kconfig2Logic>

<sup>23</sup> <https://github.com/davidfa71/Extending-Logic/tree/master/code/Logic2BDD>

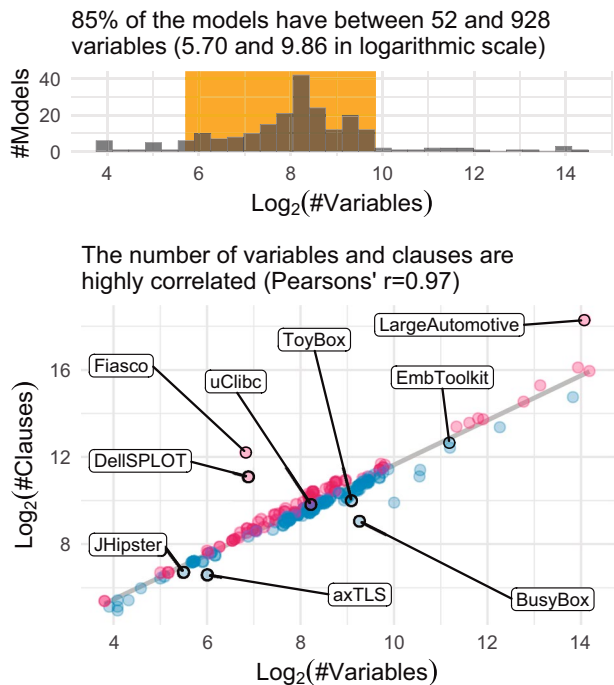
**Table 2** Software configuration models included in the benchmark

Model	#Variables	#Clauses	#SAT-Solutions
JHipster (Halin et al. 2019)	45	104	26,256
axTLS 1.5.3 ( <a href="http://axtls.sourceforge.net/">http://axtls.sourceforge.net/</a> )	64	96	$3.924 \cdot 10^{12}$
Fiasco 2014092821 ( <a href="https://os.inf.tu-dresden.de/fiasco/">https://os.inf.tu-dresden.de/fiasco/</a> )	113	4,717	$5.144 \cdot 10^9$
DellSPLOT (Nöhner and Egyed 2013)	118	2,181	$7.440 \cdot 10^6$
uClibc 201 50420 ( <a href="https://www.uclibc.org/">https://www.uclibc.org/</a> )	298	903	$7.503 \cdot 10^{50}$
ToyBox 0.5.2 ( <a href="http://landley.net/toybox/">http://landley.net/toybox/</a> )	544	1,020	$1.450 \cdot 10^{17}$
BusyBox 1.23.2 ( <a href="https://busybox.net/">https://busybox.net/</a> )	613	530	$7.428 \cdot 10^{146}$
EmbToolkit 1.7.0 ( <a href="https://www.embt toolkit.org/">https://www.embt toolkit.org/</a> )	2,331	6,437	$3.961 \cdot 10^{334}$
LargeAutomotive (Krieter et al. 2018)	17,365	321,897	$5.260 \cdot 10^{1,441}$

## 5.2 Q1: Scalability of samplers

The following experiment E1 was undertaken to obtain a sample set S1 for answering Q1. Each sampler generated a sample with one thousand configurations for every model in the benchmark. The timeout for each sample generation was set to one hour. Table 3

**Fig. 10** Size of the benchmark models in terms of the number of variables and clauses



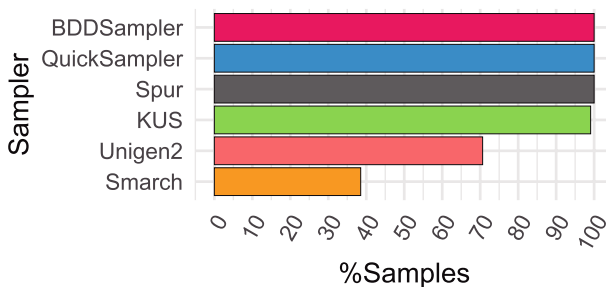
**Table 3** Sample generation time in seconds for the configuration models (sample size = 1,000 configurations; timeout = 1 hour)

Model	BDD Sampler	KUS	Quick Sampler	Smarch	Spur	Unigen2
JHipster	0.04	0.27	0.07	911.08	0.03	3.59
axTLS	0.04	0.34	0.20	1,993.90	0.03	timeout
Fiasco	0.07	0.45	1.47	timeout	0.06	timeout
DellSPLOT	0.08	0.44	0.44	3,278.09	0.07	187.58
uClibc	0.14	0.99	0.50	timeout	0.23	timeout
ToyBox	0.25	1.25	0.78	timeout	0.09	timeout
BusyBox	0.26	1.87	0.67	timeout	0.17	timeout
EmbToolkit	2.61	timeout	4.62	timeout	9.15	timeout
LargeAutomotive	12.07	119.26	77.06	timeout	24.57	timeout

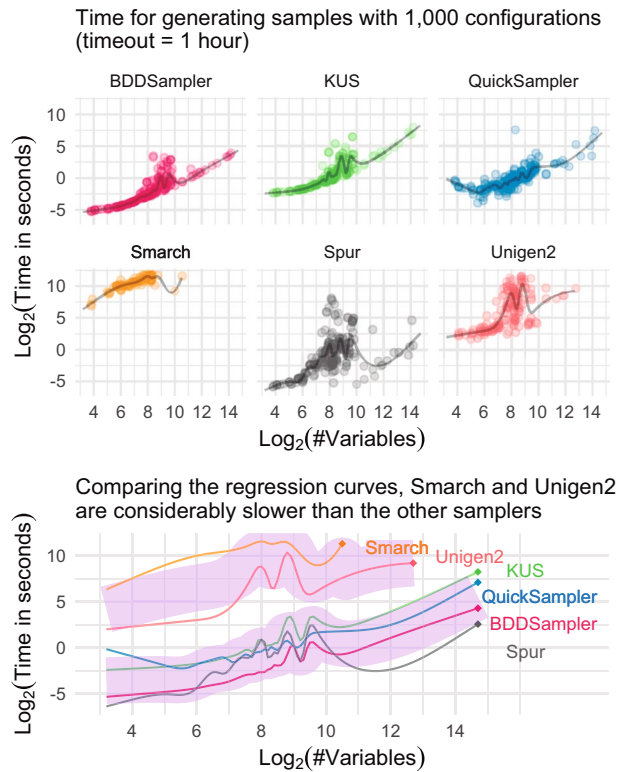
summarizes the generation times for the configuration models. The histogram in Fig. 11 shows the percentage of samples that each sampler was able to generate. In total, 257.92 hours (10.75 days) of CPU time were needed for generating the samples (or reaching the timeout).

### 5.3 Q2: Uniformity of samplers

The following experiment E2 was carried out to obtain a sample set S2 for answering Q2, Q3, and Q5. Each sampler was run to generate a sample for every model in the benchmark. As Section 5.4 will explain in detail, the number of configurations per sample was estimated for  $\alpha = 0.01$ ,  $\beta = 0.01$ , and  $w = 0.1$ . The timeout for each sample generation was set to one hour. In total, 373.5 hours (15.56 days) of CPU time were needed for generating the samples (or reaching the timeout). The histogram in Fig. 13 summarizes the results. Nearly all samples produced by BDDSampler, Smarch, Spur, and Unigen2 obtained high  $p$ -values in the range (0.9, 1]. In contrast, KUS and Quick-Sampler generated many samples with  $p$ -values in the interval [0, 0.1]. Since  $\alpha$  is set to 0.01, remember from Section 4.2 that a  $p$ -value less or equal to 0.01 means rejecting the uniformity hypothesis. Likewise, a  $p$ -value close to 1 reflects that the sample greatly supports the uniformity hypothesis. Table 4 summarizes the  $p$ -values for the configuration models in detail.

**Fig. 11** Percentage of samples that each sampler was able to generate (sample size = 1,000 configurations; timeout = 1 hour)

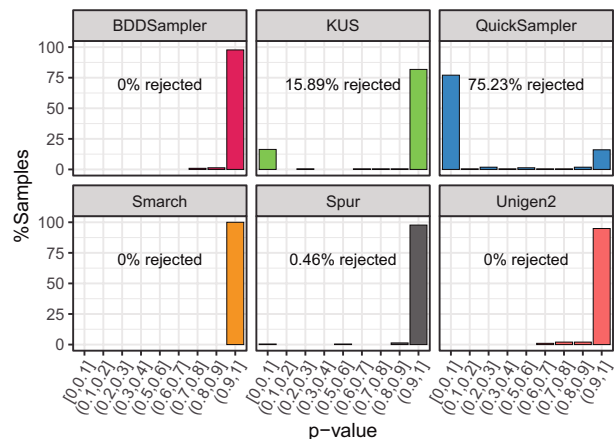
**Fig. 12** Time the samplers needed to generate 1,000 configurations for each model in the benchmark



KUS and Spur implement Knuth's sampling procedure (see Section 2.1.3). Accordingly, they should be uniform "by design". Moreover, the KUS and Spur empirical validations in Sharma et al. (2018) and Achlioptas et al. (2018), respectively, did not detect any problem (though only small models with a few hundred variables were used). However, our inspection using more varied and larger models revealed the following uniformity flaws:

- As Fig. 13 shows, 16.4% of the KUS samples got a  $p$ -value in  $[0, 0.1]$ . Furthermore, in 15.89% of the cases, the  $p$ -values were less than  $\alpha = 0.01$ , and thus rejected the uni-

**Fig. 13** Goodness-of-fit test results for the whole benchmark ( $\alpha = \beta = 0.01$  and  $w = 0.1$ )

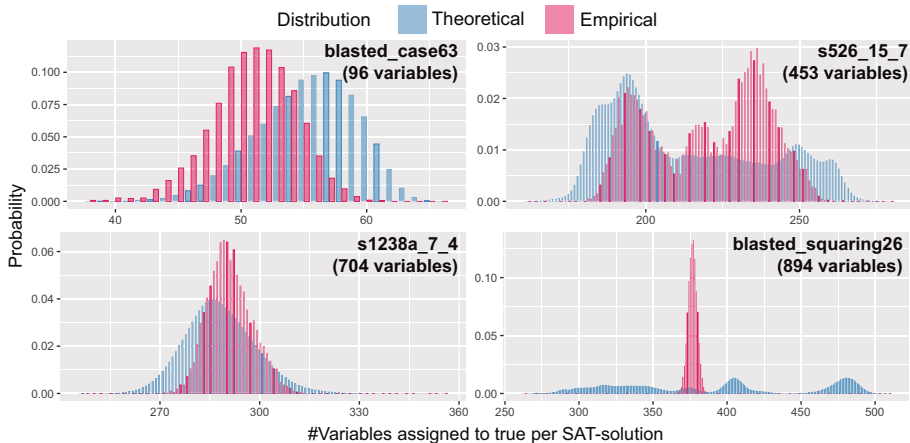


**Table 4** Goodness-of-fit  $p$ -values for the configuration models ( $\alpha = \beta = 0.01$  and  $w = 0.1$ ; timeout = 1 hour)

Model	BDD Sampler	KUS	Quick Sampler	Smarch	Spur	Unigen2
JHipster	$\sim 1$	0.99	$\sim 0$	$\sim 1$	$\sim 1$	$\sim 1$
axTLS	$\sim 1$	$\sim 1$	$\sim 0$	timeout	$\sim 1$	timeout
Fiasco	$\sim 1$	$\sim 1$	0.30	timeout	$\sim 1$	timeout
DellSPLOT	$\sim 1$	0.99	0.85	timeout	$\sim 1$	0.96
uClibc	$\sim 1$	$\sim 1$	$\sim 1$	timeout	$\sim 1$	timeout
ToyBox	$\sim 1$	$\sim 1$	$\sim 0$	timeout	$\sim 1$	timeout
BusyBox	$\sim 1$	$\sim 1$	$\sim 0$	timeout	$\sim 1$	timeout
EmbToolkit	$\sim 1$	timeout	$\sim 1$	timeout	$\sim 0$	timeout
LargeAutomotive	$\sim 1$	timeout	$\sim 0$	timeout	$\sim 1$	timeout

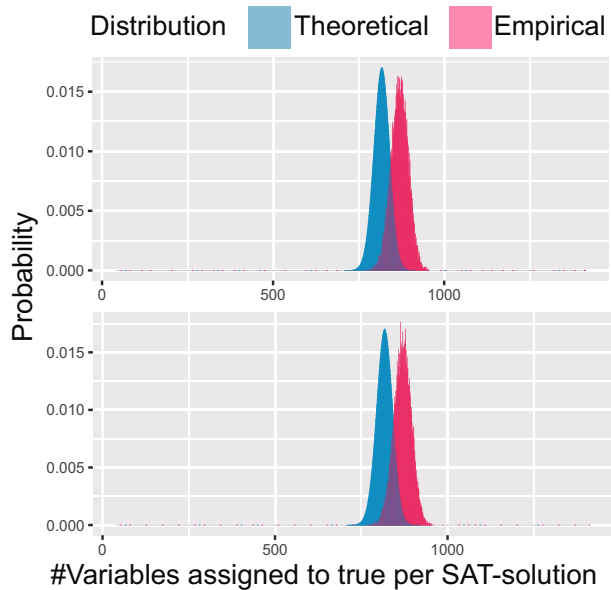
formity hypothesis. Figure 14 shows four examples where KUS uniformity was rejected. Each subfigure compares, for a particular model, the histogram of the SAT-solution distribution of the whole population (in blue) with the distribution of the generated sample (in red). Unfortunately, the rejected samples do not show any clear pattern that explains the causes of KUS failures. For instance, KUS exhibits difficulties with small models (blasted\_case63) but also with large ones (blasted\_squaring26), with normal distributions (blasted\_case63 and s1238a\_7\_4) and non-normal distributions (s526\_15\_7 and blasted\_squaring26), with left-skewed distributions (s1238a\_7\_4) and right-skewed distributions (blasted\_case63), etc.

- In our previous evaluation (Heradio et al. 2020), we detected that Spur generated uniform samples for all models except for EmbToolkit. We thought our test was making a Type 1 error, misjudging the sampler uniformity because an extremely low  $p$ -value happened due to randomness. However, when we checked the samplers' uniformity with our new test, we obtained exactly the same results for this particular model, which raised our suspicions. We repeated the experiment one thousand times and Spur never generated a uniform sample for EmbToolkit. Figure 15 shows the results for two of those experiment repetitions. In this case, Spur's error always displays the same pattern: the solutions in the sample have more variables assigned to true than in the population.


**Fig. 14** Example of KUS samples rejected with the goodness-of-fit test



**Fig. 15** Two samples Spur generated for EmbToolkit



#### 5.4 Q3: Scalability of the SFpC test

Two factors influence the scalability of a uniformity test when applied to a particular model and sampler: (i) the number of configurations the test needs to consider, and (ii) the time the test invests in analyzing those configurations.

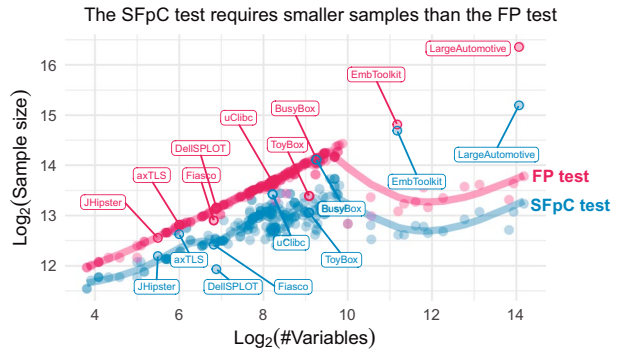
Concerning the first factor, and as discussed in Section 2, the methods proposed in the literature to verify samplers' uniformity require colossal sample sizes with millions of configurations. Thus uniformity had been tested on trivial models so far, with a few hundred variables. To support evaluating uniformity over more complex models, in Heradio et al. (2020) we proposed the FP test, which compares the variable frequency distribution of a sample with the variable probability distribution of the entire population. With this test, we could validate samplers' uniformity on models with more than seventeen thousand variables (Heradio et al. 2020). Figure 16 compares the sample sizes that the FP test needs (in red) with the sample sizes our new SFpC test requires (in blue), showing that the latter needs fewer configurations in most cases.

In Fig. 16, each model's sample size was determined with the procedure described in Section 4.2. In particular, the R package `pwr`<sup>24</sup> (Kabacoff 2011) was used to perform Cohen's power tables calculations. To ensure the highest reliability of the samplers' uniformity tests (see Section 5.4), we set  $\alpha = 0.01$ ,  $\beta = 0.01$ , and  $w = 0.1$ . That is, the  $\chi^2$  test confidence level was fixed to 99%, the power to 99%, and the effect size to *small*. Table 5 compares in detail the samples sizes obtained for the configuration models.

The sample size depends on the model's degrees of freedom in both the FP and the SFpC tests. Nevertheless, each test defines *degrees of freedom* in a different way. The

<sup>24</sup> <https://cran.r-project.org/web/packages/pwr>

**Fig. 16** Comparison of the sample sizes consumed by the FP and SFpC tests ( $\alpha = \beta = 0.01$  and  $w = 0.1$ )



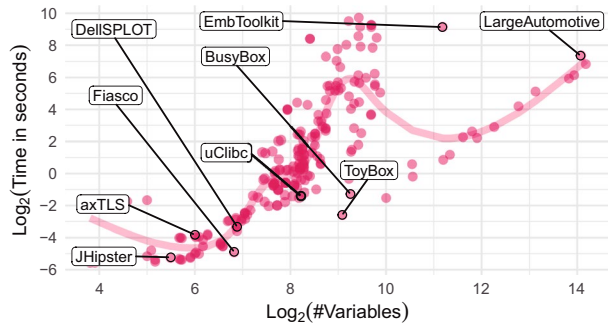
degrees of freedom of the FP test  $df_{FP}$  are the number of variables (minus one) whose probability is neither zero nor one (see Section 3 of Heradio et al. (2020)). The degrees of freedom of the SFpC test  $df_{SFpC}$  are the number cases (minus one) for which there is at least one valid configuration with a particular number of variables assigned to true (see Section 4.1). As Fig. 16 shows, in practice,  $df_{SFpC} \leq df_{FP}$  and therefore the SFpC test consumes fewer configurations.

Regarding the time SFpC requires to analyze the generated configurations, once the theoretical distribution of SAT solutions is known, the remaining computations can be performed extremely fast (see Section 4.1). So, the SFpC's potential bottleneck is getting such distribution with the algorithm PD. Figure 17 shows the time it took to compute the theoretical distribution for each model in the benchmark, ranging from 0.02 seconds to 14.14 minutes. Table 6 details the times for the configuration models. It is worth noting that the model which needed the longest time was `s1196a_3_2`, which is an industrial SAT formula (thus not included in Table 2). This illustrates the dependency that BDDs have on variable ordering heuristics. Whereas this model has a medium-size CNF formula (690 variables and 1,805 clauses), the BDD we synthesized was huge (2,284,697 nodes). In contrast, for `LargeAutomotive` (17,365 variables and 321,897) a more reduced BDD was obtained (30,432 nodes), and hence computing its theoretical SAT-solution distribution just took 2.74 minutes.

**Table 5** Sample sizes the SFpC and FP tests required for the configuration models ( $\alpha = \beta = 0.01$  and  $w = 0.1$ )

Model	Sample size for the SFpC test	Sample size for the FP test
JHipster	4,664	5,994
axTLS	6,314	7,198
Fiasco	5,460	7,646
DellSPLIT	3,889	9,131
uClibc	10,987	13,047
ToyBox	8,517	10,739
BusyBox	17,738	18,041
EmbToolkit	26,482	28,866
LargeAutomotive	37,626	84,522

**Fig. 17** Time it took to compute the distribution of SAT-solutions for all models in the benchmark



## 5.5 Q4: Validity of SFpC

Two criteria are typically used for assessing measurement quality (Trochim et al. 2015): *validity* and *reliability*. Since we are interested in the quality of SFpC measurements, *validity* will refer to what extent SFpC actually measures uniformity, and *reliability* will refer to repeatability, i.e., to the consistency of the results obtained when SFpC is applied several times to the same sampler and model. This section examines SFpC’s validity, and the next section deals with SFpC’s reliability.

To evaluate SFpC’s validity, we followed a *convergent strategy* (Trochim et al. 2015) by examining the degree to which SFpC results are similar to those obtained by other uniformity tests. Table 7 summarizes the uniformity verdicts reported in the literature. There is a total consensus that Unigen2 is uniform and QuickSampler is not. SFpC results are consistent with this consensus.

As we mentioned in Section 5.4, before the publication of FP in Heradio et al. (2019), the literature relied on limited tests that only could handle the simplest models with a few hundred variables. As more complex are considered, the chances to detect samplers’ additional flaws increases. In other words, the sensitivity of FP and SFpC is higher than their predecessors. Accordingly, we performed a new Experiment E3 focused on checking the convergent validity of FP and SFpC in detail. A new sample set S3 was procured by asking each sampler to generate a sample for every model in the benchmark. Then, the uniformity of the samples was analyzed with both FP and SFpC. Since FP generally needs larger samples than SFpC (see Section 5.4), the sample sizes were set according to FP requirements.

**Table 6** Seconds it took to compute the distribution of SAT-solutions for the configuration models

Model	Time
JHipster	0.03
axTLS	0.07
Fiasco	0.03
DellSPLOT	0.10
uClibc	0.38
ToyBox	0.16
BusyBox	0.41
EmbToolkit	567.93
LargeAutomotive	164.35

**Table 7** Samplers' uniformity judgments reported in the literature

Article	Uniform	Non-Uniform
Achlioptas et al. (2018)	Spur and Unigen2	–
Chakraborty et al. (2015)	Unigen2	–
Chakraborty et al. (2019)	Unigen2	QuickSampler
Oh et al. (2019)	Smarch and Unigen2	–
Plazar et al. (2019)	Unigen	QuickSampler
Sharma et al. (2018)	KUS and Spur	–
This present paper (SFpC)	BDDSampler, Smarch and Unigen2	KUS, QuickSampler and Spur

Due to its higher sensitivity compared to prior tests, SFpC detects that KUS and Spur sometimes behave non-uniformly

Pearson's correlation coefficient  $\rho$  of the  $p$ -values obtained with FP and SFpC was  $\rho = 0.953$ , and Cohen's kappa  $\kappa$  of the test verdicts (i.e., rejection/acceptance of sampler's uniformity) was  $\kappa = 0.942$ . As FP and SFpC results were numerically highly correlated, and their final judgments were remarkably consistent, convergent validity was successfully confirmed.

## 5.6 Q5: Reliability of SFpC

SFpC's reliability was evaluated with a *test-retest* strategy (Trochim et al. 2015) by comparing its results with the sample sets S2 and S3. Pearson's correlation coefficient of the  $p$ -values calculated with SFpC in S2 and S3 was  $\rho = 0.950$ , and Cohen's kappa of the corresponding final judgments (i.e., rejection/acceptance of sampler's uniformity) was  $\kappa = 0.939$ . As a result, SFpC's reliability was positively evaluated.

## 6 Discussion

The experimental results indicate that SFpC supports testing samplers' uniformity on complex models with thousands of variables and constraints, providing valid and reliable judgments. The results show that the only sampler that satisfies both scalability and uniformity is BDDSampler. The following points summarize the **key findings** per research question:

**Q1: Samplers' scalability.** BDDSampler, KUS, QuickSampler, and Spur are by far faster than Smarch and Unigen2. This finding agrees with the prior evaluations reported by Plazar et al. (2019) and Heradio et al. (2020).

**Q2: Samplers' uniformity.** Three categories of samplers can be distinguished: (i) those that mostly fail to produce uniform samples (QuickSampler), (ii) those that usually work but from time to time generate non-uniform samples (KUS and Spur), and (iii) those that always produce uniform samples (BDDSampler, Smarch, and Unigen2). QuickSampler's incapacity to generate uniform samples was previously reported by Chakraborty et al. (2019), Plazar et al. (2019), and Heradio et al. (2020). However, this

paper is the first one that detects problems with **KUS** and **Spur**. We think this finding is due to SFpC's ability to test samplers' uniformity on considerably more complex models than previous tests.

**Q3: SFpC's scalability.** SFpC is the most scalable uniformity test to date. It requires the smallest sample size of all existing tests, enabling the verification of samplers' uniformity in large models even for the most strict quality settings ( $\alpha = 0.01$ ,  $\beta = 0.01$ , and  $w = 0.1$ ).

**Q4: SFpC's validity.** According to the results, SFpC judgments are consistent with the verdicts given by the alternative methods proposed in the literature.

**Q5: SFpC's reliability.** The results show that SFpC judgments are reliable, i.e., when SFpC is applied repeatedly to the same model and sampler, the reached conclusions are notably consistent.

The **implications** of our research are twofold:

1. As uniform random sampling is a strong requirement for many relevant analyses on configurable systems, **BDDSampler**'s positive impact may be considerable, e.g., to test SPLs (Halin et al. 2019; Plazar et al. 2019), to support predicting and optimizing the performance of configurable systems (Oh et al. 2017; Kaltenecker et al. 2020), etc. As an illustrative example of the importance that sampling has to SPL practitioners, in the SPLC 23<sup>rd</sup> edition, there was a *challenge* dedicated specifically to this topic and entitled "Product Sampling for Product Lines: The Scalability Challenge" (Pett et al. 2019). Moreover, different papers have been recently published on uniform random sampling, and other sorts of sampling such as *t*-wise, in SPLC (Varshosaz et al. 2018; Oh et al. 2019; Muñoz et al. 2019) and the International Working Conference on Variability Modelling of Software-Intensive Systems (VaMoS) (Krieter et al. 2020). Furthermore, the applicability of **BDDSampler** goes beyond the SPL domain since sampling is also needed in artificial intelligence (Chakraborty and Meel 2019; Dutra et al. 2018; Roy et al. 2018), integrated circuit simulation and verification (Hébner and Becker 2011; Naveh et al. 2006; Yuan et al. 2002), etc.
2. SFpC can be used to debug and thus improve existing samplers (see Figs. 14 and 15), or to validate future samplers. The importance of samplers' validation is well recognized by the SPL community. Recently, in the SPLC 25<sup>th</sup> edition, there was a session dedicated to "Sampling, variability analysis and visualization", where two tools for samplers' evaluation were presented: **BURST** (Acher et al. 2021) and **AutoSMP** (Pett et al. 2021). Those tools could be enhanced by integrating SFpC; e.g., **BURST** relies on **Barbarik**, which has inferior performance than SFpC (see Section 2.2.2). Again, the interest in samplers' validation is not restricted to SPLs. In fact, most uniformity tests have been proposed by artificial intelligence researchers, mainly from the SAT community (Achlioptas et al. 2018; Chakraborty 2015; Dutra et al. 2018; Sharma et al. 2018; Chakraborty and Meel 2019).

It is worth noting that our work has the following **limitation**: both **BDDSampler** and SFpC rely on BDD technology. Synthesizing the BDD encoding of a variability model is sometimes unattainable. This is because the variable ordering chosen to build a BDD dramatically impacts its size, and finding the optimal ordering is an NP-problem. So the search is approached heuristically without guarantees. This problem principally affects **BDDSampler**, but not much SFpC.

1. **BDDSampler** receives a model's BDD encoding as input. If the available heuristics fail to find an adequate variable ordering, **BDDSampler** becomes useless, and an alternative technology (e.g., SAT) must be used.
2. SFpC can evaluate any sampler, independently of the technology in which it is built. Indeed, Section 5 reports the SFpC use for samplers implemented with BDDs (**BDDSampler**), #SAT-solvers (**QuickSampler**, **Smarch**, **Spur**, and **Unigen2**), and d-DNNFs (**KUS**). The impossibility of creating a BDD for a particular model only prevents SFpC from using it as part of the benchmark presented in Section 5.1. Currently, the benchmark includes 218 models with their respective BDDs. In our opinion, the models' great variety in terms of size (from 14 to 18,570 variables) and application domain (automotive industry, embedded systems, a laptop customization system, a web application generator, integrated circuits, etc.) is adequate to ensure samplers' verification to a great degree.

Finally, the following **threats to our study's validity** should be taken into account:

1. There is no absolute guarantee that the samplers we have certified as uniform behave non-uniformly in models not included in the benchmark.
2. Our experimental design discards two potential confounders for evaluating the scalability of samplers:
  - *Sampling parallelization.* Although any sampler can be run in a multi-core fashion, thus producing samples concurrently, only **Unigen2** and **Smarch** were specifically designed for that. The focus of our evaluation is on the sampling techniques, not on how those techniques can be parallelized efficiently. Therefore, all samplers were run on a single thread.
  - *Use of preprocessing techniques.* There are some methods to preprocess the model Boolean formulas for speeding up further computations. For example, Ivrrí et al. (2016) claim that sampling with the formulas' MIS produces 2-3 orders of magnitude performance improvement. Nevertheless, Plazar et al. (2019) empirical results contradict that, showing no running time difference between sampling from the whole formula or the MIS. Anyway, we decided to focus on the sampling techniques, not on how any additional preprocessing methods may impact those techniques.

## 7 Conclusions

The number of SAT solutions that configuration models encompass can be so large that most analyses cannot be performed neither examining every valid configuration, nor calling a SAT solver massively. Statistical inference opens an alternative way to address these problems by working with a tractable sample accurately predicts results for the entire space. However, the laws of statistical inference impose an indispensable requirement that samples must be collected at random, i.e., the configuration space needs to be covered uniformly.

Two major research challenges on SAT-solution random sampling have been addressed in this paper: we (i) developed a new random sampler, called **BDDSampler**, and (ii) proposed a goodness-of-fit test to verify samplers' uniformity. Our new test

requires the least sample size of all existing methods, in the literature, supporting the samplers' uniformity assessment even on colossal models and the most strict reliability arrangements. Using this test, we have undertaken the empirical evaluation of six state-of-the-art samplers, revealing that only **BDDSampler** satisfies both uniformity and scalability.

It is worth remarking that **BDDSampler** works with a BDD encoding of a configuration model as input, and synthesizing such BDD is not always feasible as it depends on finding an adequate variable order heuristically. Our work deals with this limitation by exposing uniformity bugs on two scalable samplers based on alternative technologies (**KUS** on d-DNNFs and **Spur** on #SAT), thus facilitating their fixing. Having available all these samplers would support coping with the varied difficulties that the Boolean encoding of configuration models poses (e.g., large intractable CNFs, enormous BDDs, etc.).

## 8 Material

Following *open science*'s good practices, our software artifacts are available publicly.

- **BDDSampler** is available at <https://github.com/davidfa71/BDDSampler>
- The code scripts to replicate our experimental validation (i.e., to calculate each model's sample size, run the samplers, and test the scalability/uniformity of the samplers) are available at <https://github.com/rheradio/ConfSystSampling>
- A detailed report on every research question in Section 5 is available at: <https://rheradio.github.io/ConfSystSampling>
- The data of Experiments E1 and E2 (including the benchmark models in DIMACS/DDDMP/NNF formats, the generated samples, the goodness-of-fit test results, etc.) are available at <https://doi.org/10.5281/zenodo.4514919>
- The data of Experiment E3 are available at <https://doi.org/10.5281/zenodo.5509947>

**Acknowledgements** This work has been partially funded by the Universidad Nacional de Educacion a Distancia (project OPTIVAC 096-034091 2021V/PUNED/008); the Spanish Ministry of Science, Innovation and Universities (project OPHELIA RTI2018-101204-B-C22); the Community of Madrid (research network ROBOCITY2030-DIH-CM S2018/NMT-4331); the TASOVA network (MCIU-AEI TIN2017-90644-REDT); and the Junta de Andalucía (METAMORFOSIS project).

**Funding** Open Access funding provided thanks to the CRUE-CSIC agreement with Springer Nature.

## Declarations

**Conflicts of interest** The authors have no conflict of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Acher M, Perrouin G, Cordy M (2021) BUrst: a benchmarking platform for uniform random sampling techniques. In: 25th systems and software product line conference (SPLC). Leicester, United Kingdom
- Achlioptas D, Hammoudeh ZS, Theodoropoulos P (2018) Fast sampling of perfectly uniform satisfying assignments. In: 21st International conference on theory and applications of satisfiability testing (SAT). Oxford, UK, pp 135–147
- Alf  rez M, Acher M, Galindo J, Baudry B, Benavides D (2019) Modeling variability in the video domain: language and experience report. *Software Quality Journal* 27(1):307–347
- Alves Pereira J, Acher M, Martin H, J  z  quel JM (2020) Sampling effect on performance prediction of configurable systems: A case study. In: ACM/SPEC international conference on performance engineering (ICPE). Edmonton AB, Canada, pp 277–288
- Batory DS (2005) Feature models, grammars, and propositional formulas. In: 9th software product line conference (SPLC). Rennes, France, pp 7–20
- Bellare M, Goldreich O, Petrank E (2000) Uniform generation of NP-witnesses using an NP-oracle. *Information and Computation* 163(2):510–526
- Berger T, She S, Lotufo R, Wasowski A, Czarnecki K (2013) A study of variability models and languages in the systems software domain. *IEEE Transactions on Software Engineering* 39(12):1611–1640
- Biere A, Heule M, van Maaren H, Walsh T (2009) Handbook of satisfiability: Volume 185 frontiers in artificial intelligence and applications. IOS Press, Amsterdam
- Brace KS, Rudell RL, Bryant E (1990) Variability modeling in the real: a perspective from the operating systems domain. In: IEEE/ACM design automation conference (DAC). Orlando, Florida, USA, pp 40–45
- Bryant RE (1986) Graph-based algorithms for boolean function manipulation. *IEEE Transactions on Computers* C-35(8):677–691
- Chakraborty S, Fremont DJ, Meel KS, Seshia SA, Vardi MY (2015) On parallel scalable uniform SAT witness generation. In: 21st international conference on tools and algorithms for the construction and analysis of systems (TACAS). London, UK, pp 304–319
- Chakraborty S, Meel KS (2019) On testing of uniform samplers. In: 33rd conference on artificial intelligence (AAAI). Honolulu, Hawaii, USA, pp 7777–7784
- Chakraborty S, Meel KS, Vardi MY (2013) A scalable and nearly uniform generator of SAT witnesses. In: 25th international conference on computer aided verification (CAV). Saint Petersburg, Russia, pp 608–623
- Chihara LM, Hesterberg TC (2011) Mathematical statistics with Resampling and R. Wiley, New York
- Chollet F, Allaire J (2018) Deep learning with R. Manning Publications
- Cohen J (1988) Statistical power analysis for the behavioral sciences. Routledge, Evanston
- Cover TM, Thomas JA (2006) Elements of information theory. Wiley, New York
- D’Agostino RB, Stephens MA (1986) Goodness-of-fit-techniques. CRC Press, Boca Raton
- Davis M, Logemann G, Loveland D (1962) A machine program for Theorem-proving. *Communications of the ACM* 5(7):394–397
- van Dijk, T (2016) Sylvan: multi-core decision diagrams. Ph.D. thesis, University of Twente
- Dutra R, Laeufer K, Bachrach J, Sen K (2018) Efficient sampling of SAT solutions for testing. In: 40th international conference on software engineering (ICSE). ACM, New York, NY, USA, pp 549–559
- Faul F, Erdfelder E, Lang AG, Buchner A (2007) G\*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods* 2(39):175–191
- Fernandez-Amoros D, Bra S, Aranda-Escobal   E, Heradio R (2020) Using extended logical primitives for efficient BDD building. *Mathematics* 8(8):1–17
- Fernandez-Amoros D, Heradio R, Mayr-Dorn C, Egyed A (2019) A Kconfig translation to logic with one-way validation system. In: 23rd international systems and software product line conference (SPLC). Paris, France, pp 303–308
- Goodfellow I, Bengio Y, Courville A (2016) Deep learning. The MIT Press, Cambridge
- Grosse I, Bernaola-Galvan P, Carpena P, Roman-Roldan R, Oliver J, Stanley HE (2002) Analysis of symbolic sequences using the Jensen-Shannon divergence. *Physical Review E* 65(2):041,905/1-0419,051/16
- Guo J, Yang D, Siegmund N, Apel S, Sarkar A, Valov P, Czarnecki K, Wasowski A, Yu H (2018) Data-efficient performance learning for configurable systems. *Empirical Software Engineering* 23(3):1826–1867
- Halin A, Nuttinck A, Acher M, Devroey X, Perrouin G, Baudry B (2019) Test them all, is it worth it? Assessing configuration sampling on the JHipster Web development stack. *Empirical Software Engineering* 24(2):674–717



- Heradio R, Fernandez-Amoros D, Galindo JA, Benavides D (2020) Uniform and scalable SAT-sampling for configurable systems. In: 24th systems and software product line conference (SPLC). Montréal, Canada, pp 1–11
- Heradio R, Fernandez-Amoros D, Mayr-Dorn C, Egyed A (2019) Supporting the statistical analysis of variability models. In: 41st international conference on software engineering (ICSE). Montréal, Canada, pp 843–853
- Hëbner M, Becker J (2011) Multiprocessor system-on-Chip: Hardware design and tool integration. Springer, New York
- Ivrii A, Malik S, Meel KS, Vardi MY (2016) On computing minimal independent support and its applications to sampling and counting. *Constraints* 21(1):41–58
- Kabacoff R (2011) *R in Action: Data analysis and graphics with R*. Manning Publications
- Kaltenecker C, Grebhahn A, Siegmund N, Apel S (2020) The interplay of sampling and machine learning for software performance prediction. *IEEE Software* 37(4):58–66
- Kaltenecker C, Grebhahn A, Siegmund N, Guo J, Apel S (2019) Distance-based sampling of software configuration spaces. In: 41st international conference on software engineering (ICSE). Montreal, Canada, pp 1084–1094
- Kaplan D (2012) Statistical modeling: A fresh approach. Project Mosaic
- Knuth DE (2009) The art of computer programming, Volume 4, Fascicle 1: Bitwise Tricks & Techniques; Binary Decision Diagrams. Addison-Wesley Professional
- Kolesnikov S, Siegmund N, Kästner C, Grebhahn A, Apel S (2019) Tradeoffs in modeling performance of highly configurable software systems. *Software & Systems Modeling* 18(3):2265–2283
- Krieter S (2019) Enabling efficient automated configuration generation and management. In: 23rd international systems and software product line conference (SPLC). Paris, France, pp 215–221
- Krieter S, Thüm T, Schulze S, Saake G, Leich T (2020) Yasa: Yet another sampling algorithm. In: 14th international working conference on variability modelling of software-intensive system (VaMoS). Magdeburg, Germany, pp 1–10
- Krieter S, Thüm T, Schulze S, Schröter R, Saake G (2018) Propagating configuration decisions with modal implication graphs. In: 40th international conference on software engineering (ICSE). Gothenburg, Sweden, pp 898–909
- Lagniez JM, Marquis P (2017) An improved decision-DNNF compiler. In: 26th International joint conference on artificial intelligence (IJCAI). Melbourne, Australia
- Lakens D (2013) Calculating and reporting effect sizes to facilitate cumulative science: a practical primer for *t*-tests and ANOVAs. *Frontiers in Psychology* 4(863):1–12
- Lin J (1991) Divergence measures based on the shannon entropy. *IEEE Transactions on Information Theory* 37(1):145–151
- Meinel C, Theobald T (1998) Algorithms and data structures in VLSI design: OBDD - foundations and applications. Springer, New York
- Mendonça, M (2009) Efficient reasoning techniques for large scale feature models. Ph.D. thesis, University of Waterloo
- de Moura L, Bjørner N (2008) Z3: An Efficient SMT Solver. In: 14th International conference on tools and algorithms for the construction and analysis of systems (TACAS). Budapest, Hungary, pp 337–340
- Muñoz DJ, Oh J, Pinto M, Fuentes L, Batory D (2019) Uniform random sampling product configurations of feature models that have numerical features. In: 23rd international systems and software product line conference (SPLC). Paris, France, pp 289–301
- Nair V, Menzies T, Siegmund N, Apel S (2017) Using bad learners to find good configurations. In: 11th joint meeting on foundations of software engineering (ESEC/FSE). Paderborn, Germany, pp 257–267
- Narodytska N, Walsh T (2007) Constraint and variable ordering heuristics for compiling configuration problems. In: International joint conference on artificial intelligence (IJCAI). Hyderabad, India, pp 149–154
- Naveh Y, Rimón M, Jaeger I, Katz Y, Vinov M, Marcus E, Shurek G (2006) Constraint-based random stimuli generation for hardware verification. In: 18th innovative applications of artificial intelligence conference (IAAI). Boston, Massachusetts, USA, pp 1720–1727
- Nöhner A, Egyed A (2013) C2O configurator: a tool for guided decision-making. *Automated Software Engineering* 20(2):265–296
- Oh J, Batory D, Myers M, Siegmund N (2017) Finding near-optimal configurations in product lines by random sampling. In: 11th joint meeting on foundations of software engineering (ESEC/FSE). New York, NY, USA, pp 61–71
- Oh J, Gazzillo P, Batory DS (2019) t-wise coverage by uniform sampling. In: 23rd international systems and software product line conference (SPLC). Paris, France, pp 84–87

- Pett T, Krieter S, Thüm T, Lochau M, Schaefer I (2021) AutoSMP: An evaluation platform for sampling algorithms. In: 25th systems and software product line conference (SPLC). Leicester, United Kingdom
- Pett T, Thüm T, Runge T, Krieter S, Lochau M, Schaefer I (2019) Product sampling for product lines: The scalability challenge. In: 23rd international systems and software product line conference (SPLC). Paris, France, pp 78–83
- Plazar Q, Acher M, Perrouin G, Devroey X, Cordy M (2019) Uniform sampling of SAT solutions for configurable systems: Are We There Yet? In: 12th IEEE conference on software testing, validation and verification (ICST). Xian, China, pp 240–251
- Roy S, Pandey A, Dolan-Gavitt B, Hu Y (2018) Bug synthesis: Challenging bug-finding tools with deep faults. In: 26th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering (ESEC/FSE). Lake Buena Vista, Florida, USA, pp 224–234
- Sharma S, Gupta R, Roy S, Meel KS (2018) Knowledge compilation meets uniform sampling. In: 22nd international conference on logic for programming, artificial intelligence and reasoning (LPAR). Awassa, Ethiopia, pp 620–636
- van Solingen R, Berghout E (1999) The goal/question/metric method: A practical guide for quality improvement of software development. McGraw-Hill, New York
- Temple P, Galindo JA, Acher M, Jézéquel JM (2016) Using machine learning to infer constraints for product lines. In: 20th international systems and software product line conference (SPLC). Beijing, China, pp 209–218
- Thurley M (2006) sharpSAT - Counting models with advanced component caching and implicit BCP. In: 9th international conference on theory and applications of satisfiability testing (SAT). Seattle, WA, USA, pp 424–429
- Trochim WM, Donnelly JP, Arora K (2015) Research methods: the essential knowledge base. Wadsworth Publishing
- Varshosaz M, Al-Hajjaji M, Thüm T, Runge T, Mousavi MR, Schaefer I (2018) A classification of product sampling for software product lines. In: 22nd International Systems and Software Product Line Conference (SPLC). Gothenburg, Sweden, pp 1–13
- Vasishth S, Broe M (2011) The foundations of statistics: a simulation-based approach. Springer, New York
- Weckesser M, Kluge R, Pfannemüller M, Matthé M, Schürr A, Becker C (2018) Optimal reconfiguration of dynamic software product lines based on performance-influence models. In: 22nd international systems and software product line conference (SPLC). Gothenburg, Sweden, pp 98–109
- Winter B (2020) Statistics for linguists: an introduction using R. Routledge, Evanston
- Yuan J, Albin K, Aziz A, Pixley C (2002) Simplifying constraint solving in random simulation generation. In: 11th IEEE/ACM international workshop on logic & synthesis (IWLS). New Orleans, Louisiana, USA, pp 185–190