

# DivSampCA: A Tuple-oriented Adaptive Sampling for Generating Small Pairwise Covering Arrays

Anonymous Authors

**Abstract**—Combinatorial interaction testing (CIT) is an effective method for verifying highly configurable systems, with pairwise testing being the most commonly used CIT technique due to its strong defect detection capabilities. A major challenge in pairwise testing is the pairwise covering array generation (PCAG) problem, which aims to generate a minimal test suite that covers all valid pairwise option value combinations. Existing sampling-based PCAG methods often generate large test suites because they do not adequately focus on directly optimizing pairwise tuple coverage when generating test cases, and they fail to exploit opportunities to use a single test case to cover multiple tuples simultaneously. To address these limitations, we propose *DivSampCA*, which employs a tuple-oriented adaptive sampling technique to enhance the diversity of the generated test suite by directly optimizing tuple coverage. Moreover, it utilizes a full coverage strategy to maximize the number of tuples covered by a single test case, ensuring that all valid pairwise combinations are covered with as few test cases as possible. We validated *DivSampCA* on 124 publicly available PCAG benchmark instances, demonstrating that our approach generates covering arrays approximately 25% smaller than those produced by the state-of-the-art PCAG algorithms. These results indicate that *DivSampCA*, which effectively reduces the size of covering arrays while maintaining full coverage, is a significant advancement in sampling techniques for solving the PCAG problem.

**Index Terms**—Pairwise Testing, Covering Array, Adaptive Sampling

## I. INTRODUCTION

With the rapid advancement of technology and the growing diversity of business requirements, software systems increasingly demand a high degree of configurability to adapt to various application scenarios [1], [2]. A highly configurable system is characterized by numerous configuration options, enabling users to tailor different objects of the program to meet specific requirements [3]. Customers can assign different values to these options to suit their individual needs.

In practical scenarios, testing all possible configurations of a configurable system to identify interaction errors between different options is inefficient or even impractical. For small-scale systems with a small number of configuration options, exhaustive testing may still be feasible. However, as the number of configuration options increases, the potential configurations grow exponentially. For example, a system with 30 options, each having 2 possible

values, would require  $2^{30} = 1,073,741,824$  tests, which is practically impossible.

To address these challenges, *combinatorial interaction testing* (CIT) constructs a covering array by selecting values for the configuration options to generate a test suite [4]. A covering array is a collection of test cases that ensures, for a fixed value  $t$ , any valid  $t$ -wise combinations of option values are included in at least one test case, thereby enabling the detection of errors caused by  $t$ -wise interactions. Considering that achieving full coverage takes a long time for large values of  $t$ , *pairwise testing* (CIT with  $t = 2$ ) is commonly adopted in practice. Pairwise testing aims to generate a *pairwise covering array* (PCA) that includes all possible value pairs. Extensive empirical studies on real-world configurable systems have shown that most defects can be detected through pairwise testing [5]–[7], highlighting its efficacy as a testing methodology.

Existing methods for generating PCAs can be broadly categorized into exact (e.g., [8]–[13]), greedy (e.g., [14]–[22]), meta-heuristic (e.g., [23]–[29]), and sampling (e.g., [30]–[37]) approaches. Exact methods generate optimal PCAs through exhaustive search, but are only feasible for small-scale problems due to their high computational complexity. Greedy algorithms build PCAs by adding one test case at a time or parameter by parameter. Although this approach can generate PCAs quickly, it may struggle with scalability issues [38], making it difficult to generate effective PCAs for large-scale problems. Meta-heuristic techniques, which leverage global search and local optimization strategies, are currently considered the best approach for generating small PCAs, though they can suffer from long processing times [39].

In recent years, advancements in formal reasoning tools have driven extensive research into sampling methods for the efficient construction of PCAs [40]–[42]. Existing sampling-based approaches generate diverse test cases by employing different sampling strategies and then select some of them to construct test suites. However, these methods merely focus on diversity and do not directly aim to enhance tuple coverage. They usually set sampling probabilities by only considering the distribution of individual variables while ignoring the interaction information between two-dimensional options. This leads to the situation where the sampled test cases fail to cover more tuples, thus increasing the size of the test suites. On the other hand, we have observed that it is highly inefficient to cover the small number of tuples that are

not covered by the current test suite by generating a large number of test cases. This is because these tuples are difficult to be covered through sampling methods, and instead, deliberately designed test cases are required. Additionally, it has been observed that tuples can coexist within the same test case, which can help to avoid the use of redundant test cases in achieving full coverage.

To address these limitations, we propose a two-phase adaptive sampling algorithm named *DivSampCA*. During the sampling phase, *DivSampCA* introduces a new technique called *tuple-oriented adaptive sampling (TOAS)*. This technique iteratively updates the sampling probabilities based on the proportion of each literal in the uncovered tuples, thereby fostering the generation of diverse test cases that cover a greater number of tuples. During the full coverage phase, *DivSampCA* incorporates an innovative approach dubbed *more tuples per case (MTPC)*, which constructs test cases capable of covering multiple uncovered tuples simultaneously. By leveraging the synergistic interaction between these two core techniques, *DivSampCA* effectively addresses the limitations of existing sampling-based methods, providing a more efficient solution to the *pairwise covering array generation (PCAG)* problem.

We conduct extensive experiments to evaluate the performance of *DivSampCA*. Specifically, we used 124 publicly available instances modeled from real-world configurable systems as benchmarks. *DivSampCA* was compared with state-of-the-art PCAG algorithms, including *SamplingCA* [37], *Campactor* [43] and *LS-Sampling-Plus* [35]. The experimental results show that *DivSampCA* can not only generate smaller PCAs more quickly than existing full coverage algorithms but also achieve higher pairwise coverage in less time compared with high  $t$ -coverage algorithms. These results demonstrate that *DivSampCA*, through the use of *TOAS* and *MTPC*, represents a significant advancement in addressing the PCAG problem.

In summary, the contributions of this paper are as follows:

- We propose a tuple-oriented adaptive sampling technique, termed *TOAS*, which enhances diversity in test suites during the sampling process.
- We develop an efficient full covering technique called *MTPC*, designed to minimize the number of test cases required to cover all remaining uncovered tuples.
- We conduct extensive experiments to evaluate the performance of *DivSampCA*. The results show that *DivSampCA* effectively generates smaller covering arrays, marking a notable improvement in PCAG algorithm development.

The subsequent sections of this paper are organized as follows: Section II introduces the relevant fundamental concepts involved in this study. Section III introduces the *DivSampCA* algorithm, detailing its key components. Section IV describes our experimental settings, while Section V presents the experimental results and provides a

comprehensive evaluation of *DivSampCA*'s performance. Section VI briefly reviews the related studies in the field. Finally, Section VII summarizes the paper and presents the conclusions.

## II. PRELIMINARIES

### A. Boolean Formulae

Given a Boolean *variable*  $x$ , a *literal* is either a variable  $x$  or its complement  $\neg x$ . A Boolean clause is composed of literals connected by logical operators such as AND ( $\wedge$ ), OR ( $\vee$ ) and NOT ( $\neg$ ). The primary form in which a Boolean formula is structured is the *conjunctive normal form (CNF)*, where a formula is a conjunction of clauses, and each clause is a disjunction of literals.

A Boolean variable  $x$  can take two values: 0 or 1. For a Boolean formula  $F$ , an *assignment function*  $\phi$  assigns a value  $\phi_i$  to the variable  $x_i$  in  $F$ , where  $\phi_i$  is either 0 or 1. A *complete assignment* is an assignment function where every variable in  $F$  is assigned a value; otherwise, it is a *partial assignment*. A *satisfying assignment* of  $F$  is a complete assignment that makes  $F$  evaluate to true. If there exists at least one satisfying assignment of  $F$ , then  $F$  is satisfiable; otherwise, unsatisfiable.

**Example 1.** Consider a CNF formula represented as:

$$F = (x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_4) \wedge (x_2 \vee x_3 \vee x_4).$$

For the given formula, a satisfying assignment  $\phi$  is:  $\{x_1, x_2, x_3, x_4\} \rightarrow \{0, 1, 1, 0\}$ . Under this assignment, each clause in the formula is evaluated to be true, thus establishing that the formula is satisfiable.

### B. SAT Solving Techniques

It is widely recognized that highly configurable software can be modeled as Boolean formulas [44]. The *Boolean satisfiability (SAT)* problem aims to determine whether there exists a truth assignment to make a Boolean formula evaluate to true. Therefore, the technology for obtaining satisfying assignments for CNF formulas is crucial for testing highly configurable software.

*MiniSAT* [45] is a lightweight and efficient SAT solver that has been widely used in the field of SAT solving research. It employs the *DPLL* mechanism to perform logical reasoning [46]. *MiniSAT* iteratively assigns values to unassigned variables and learns new clauses from conflicts to prevent the solver from revisiting the same unsatisfiable states. This process continues until all variables are assigned or a contradictory state is reached.

*ContextSAT* [37] is a SAT algorithm based on *MiniSAT*, designed to generate test cases for configurable systems. It takes a CNF formula  $F$ , a (partial) assignment  $\alpha$ , a reference assignment  $\gamma$  of variables in  $F$  (denoted as  $V(F)$ ), and a variable order  $\pi$  as inputs. The output is a solution to formula  $F$  that is close to  $\gamma$ . *ContextSAT* specifies a basic assignment and seeks solutions similar to that assignment, which helps increase the diversity of test cases. Further implementation details can be found in [37].

### C. Pairwise Covering Array Generation

A *covering array*, denoted as  $CA(N; t, m, v)$ , is an  $N \times m$  array. The parameter  $t$  represents the *strength* of coverage,  $m$  denotes the number of options, and  $v$  signifies the number of values associated with each option. In all subsequent parts of this paper,  $v$  is equal to 2. An  $N \times t$  subarray is a matrix that results from choosing  $t$  columns out of the  $m$  columns, with each row representing a  $t$ -wise tuple. The covering array itself ensures that every such  $N \times t$  subarray encompasses all possible  $t$ -wise tuples of the  $v$  values.

In combinatorial interaction testing, a test suite is typically formatted as a covering array. Each row of the covering array represents a distinct test case. Meanwhile, each column represents an option, and the elements within a column denote the values chosen for that option. When  $t$  is set to 2, the array is known as a pairwise covering array, and the PCAG problem refers to the construction of an optimal PCA that meets specific constraints.

### III. OUR PROPOSED APPROACH

In this section, we describe the implementation of *DivSampCA*. Inspired by the two-stage framework [37], [41], we integrate tuple-oriented adaptive sampling to generate a diverse set of test cases and cover as many tuples as possible within each case for efficient full coverage. Specifically, we first introduce the overall framework of *DivSampCA*, followed by a detailed description of the core algorithmic techniques.

#### A. Overall Framework of *DivSampCA*

The overall framework of the *DivSampCA* algorithm is detailed in Algorithm 1. It takes a CNF formula  $F$ , a partial assignment  $\alpha$ , and hyperparameters  $k$  and *count* as input, and outputs the corresponding PCA if  $F$  is satisfiable. Before starting any operation, *DivSampCA* utilizes a preprocessor named *Coprocessor* [47] to simplify the CNF formula. The preprocessor possesses the property of equivalence preservation, which can enhance the solving process's efficiency without compromising correctness. The  $F$  in Algorithm 1 actually represents the simplified CNF formula. A valid test case is essentially a solution to the formula  $F$ . In lines 1-3 of the algorithm, *DivSampCA* invokes a *DPLL* algorithm to obtain the initial test case and initializes the test suite with it. One can specify the tuples to be included in the initial test case by designating the partial assignment  $\alpha$ . Subsequently, lines 4-9 of the algorithm represent the sampling phase. *DivSampCA* employs the *TOAS* algorithm to sample  $k$  (a hyperparameter) candidate test cases and stores them in  $C = \{c_1, c_2, \dots, c_k\}$ . Following this, it adds the test case that maximizes the number of newly added valid tuples to the test suite. Concretely, the *gain* function in line 6 calculates the number of tuples newly covered by each test case when it is added to the test suite. Then, in line 7, if the number of new tuples added falls short of the hyperparameter *count*, the sampling phase is terminated, and the algorithm

---

#### Algorithm 1: Overall Framework of *DivSampCA*

---

**Input** :  $F$ : Boolean formula in CNF;  
 $\alpha$  : (partial) assignment of  $F$ ;  
 $k$  : size of the candidate test case set;  
*count* : threshold for sampling benefit;

**Output**:  $T$ : pairwise covering array of  $F$ ;

```

1 if DPLL( $F, \alpha$ ) reports "No Solution" then return
   $\emptyset$ ;
2  $\alpha \leftarrow$  DPLL( $F, \alpha$ );
3  $T \leftarrow \{\alpha\}$ ;
4 while True do
  // Sample a test case set of size k
  // using the TOAS method
5    $C \leftarrow$  TOAS( $F, k$ );
  // gain( $c_i, T$ ): number of valid tuples
  // newly covered by the test case  $c_i$ 
6    $\beta \leftarrow \arg \max_{c_i} \text{gain}(c_i, T), 1 \leq i \leq k$ ;
7   if gain( $\beta, T$ ) < count then break;
8    $T \leftarrow T \cup \{\beta\}$ ;
9 end
10 while number of uncovered tuples > 0 do
  // Construct a test case that covers as
  // many tuples as possible
11    $\delta \leftarrow$  MTPC( $F$ );
12    $T \leftarrow T \cup \{\delta\}$ ;
13 end
14 return  $T$ ;
```

---

transitions into the full coverage phase. Finally, lines 10-14 of the algorithm represent the full coverage phase. *DivSampCA* iteratively invokes the *MTPC* algorithm to construct test cases that cover as many tuples as possible until completing the construction of the PCA. In the actual implementation of the algorithm, different values of the hyperparameters  $k$  and *count* influence both the size of the final PCA and the time consumed. Their impact on algorithm performance is experimentally evaluated in Section V-D. The implementations of the *TOAS* and *MTPC* will be detailed in the following two subsections.

#### B. Tuple-oriented Adaptive Sampling

During the sampling phase, *DivSampCA* adheres to a greedy framework, aiming to identify the test case that can cover the maximum number of uncovered tuples in each iteration. However, it is unlikely that an efficient algorithm (polynomial-time algorithm) for identifying such cases exists [48]. Therefore, an alternative approach is to generate some test cases that are significantly different from the existing test suite, and then select the test case that can add the maximum number of tuples from them and incorporate it into the test suite.

In practical implementation, *DivSampCA* employs a weight function to dynamically update the sampling probabilities of literals. This enables the sampling process to

---

**Algorithm 2:** *TOAS* ( $F, k$ )

---

**Input** :  $F$  : Boolean formula in CNF;  
 $k$  : size of the candidate test case set;  
**Output:**  $C$  : a set of test cases of size  $k$ ;

```
1  $C \leftarrow \emptyset$ ;  
2 for  $i \leftarrow 1$  to  $k$  do  
3   foreach  $x$  in  $V(F)$  do  
4      $P(x) \leftarrow 0$ ; //  $P(x)$ : Probability of  $x$   
       being assigned 1  
5     foreach  $\sigma$  in  $[1, 0]$  do  
6        $W(x, \sigma) \leftarrow 0$ ;  
7       foreach  $y$  in  $V(F)$  and  $y \neq x$  do  
8          $u[x, \sigma] \leftarrow$  number of uncovered tuples  
           for  $(y, x)$  with  $x = \sigma$ ;  
9          $\omega[x] \leftarrow$  number of uncovered tuples  
           for  $(y, x)$ ;  
10         $W(x, \sigma) \leftarrow W(x, \sigma) + u[x, \sigma]/\omega[x]$ ;  
11      end  
12    end  
13     $P(x) \leftarrow W(x, 1)/(W(x, 1) + W(x, 0))$   
14  end  
15   $\gamma_i \leftarrow$  sample a assignment according to  $P(x)$ ;  
16   $\pi_i \leftarrow$  a random variable order;  
17   $\beta_i \leftarrow \text{ContextSAT}(F, \emptyset, \gamma_i, \pi_i)$ ;  
18   $C \leftarrow C \cup \{\beta_i\}$ ;  
19 end  
20 return  $C$ ;
```

---

more comprehensively explore the areas that have not been adequately covered, ultimately generating diverse test cases. The intuitive idea behind this is to decrease the sampling probabilities of the literals corresponding to the tuples that have been widely covered. Meanwhile, it increases the probability of sampling the literals corresponding to the uncovered tuples. As a result, the generated test cases are less likely to overlap with the existing test suite, thereby achieving the goal of covering more uncovered tuples.

We introduce  $W(\cdot)$  as a weight function for a propositional formula  $F$ . It maps a variable assignment in  $F$  to a non-negative number. Specifically, the weight of variable  $x$  being assigned the value  $\sigma$  is calculated as follows:

$$W(x, \sigma) = \sum_{y \in V(F), y \neq x} \frac{\#uncovered(x, y, \sigma)}{\#uncovered(x, y)}, \quad (1)$$

where  $V(F)$  is the set of variables in  $F$ , and  $y$  represents a variable whose potential identities are every variable in  $F$  except for  $x$ . The notation  $\#uncovered(x, y)$  denotes the number of valid uncovered tuples corresponding to variables  $x$  and  $y$ ;  $\#uncovered(x, y, \sigma)$  signifies the number of tuples within these valid uncovered tuples in which the variable  $x$  is assigned the value  $\sigma$ . In each iteration, *DivSampCA* keeps track of which tuples remain uncovered

and updates the weight of each literal to calculate the probability of variable assignments in the next iteration. Precisely, the probability  $P(x)$  of variable  $x$  being assigned 1 is calculated as follows:

$$P(x) = \frac{W(x, 1)}{W(x, 1) + W(x, 0)}. \quad (2)$$

In the implementation of *DivSampCA*,  $\#uncovered(x, y)$  is calculated as 4 (i.e., the possible number of pairwise tuples) minus the number of tuples  $(x, y)$  already covered, regardless of whether the remaining uncovered tuples are valid. This is due to the fact that verifying the validity of each tuple before sampling would require a significant amount of time. While a more precise count might intuitively result in superior sampling outcomes, it would also degrade the overall performance of the algorithm.

**Example 2.** Consider a scenario involving four variables  $x_1, x_2, x_3$ , and  $x_4$ . The uncovered tuples are as follows:

$$\begin{bmatrix} [x_1 = 0, x_2 = 1], \\ [x_1 = 1, x_2 = 0], \\ [x_1 = 0, x_3 = 0], \\ [x_1 = 0, x_4 = 1], \\ [x_1 = 1, x_4 = 0], \\ [x_2 = 1, x_3 = 1] \end{bmatrix}$$

To determine the probability of each variable being assigned 1 in the next iteration, the weight is calculated as follows:

$$\begin{aligned} W(x_1, 1) &= \frac{\#uncovered(x_1, x_2, 1)}{\#uncovered(x_1, x_2)} + \frac{\#uncovered(x_1, x_3, 1)}{\#uncovered(x_1, x_3)} + \\ &\quad \frac{\#uncovered(x_1, x_4, 1)}{\#uncovered(x_1, x_4)} \\ &= 1/2 + 0 + 1/2 = 1 \end{aligned}$$

$\#uncovered(x_1, x_2)$  represents the number of uncovered tuples  $(x_1, x_2)$ , which equals 2 (namely,  $\{x_1 = 0, x_2 = 1\}$  and  $\{x_1 = 1, x_2 = 0\}$ ).  $\#uncovered(x_1, x_2, 1)$  represents the number of uncovered tuples  $(x_1, x_2)$  where  $x_1$  is assigned the value 1, which equals 1 (i.e.,  $\{x_1 = 1, x_2 = 0\}$ ). The calculations for other cases are similar.

Thus, the probability of  $x_1$  being assigned 1 is:

$$P(x_1) = \frac{W(x_1, 1)}{W(x_1, 1) + W(x_1, 0)} = \frac{1}{1+2} \approx 0.33$$

Similarly, the probabilities of  $x_2, x_3$  and  $x_4$  being assigned 1 are:

$$\begin{aligned} P(x_2) &= \frac{3/2}{3/2+1/2} = 0.75, P(x_3) = \frac{1}{1+1} = 0.5, P(x_4) = \\ &\quad \frac{1/2}{1/2+1/2} = 0.5 \end{aligned}$$

The specific implementation of the *TOAS* algorithm is presented in Algorithm 2. Concretely, we maintain a data structure to track the uncovered tuples in each iteration. As new test cases are selected, this data structure is updated accordingly to reflect the current coverage status. In lines 5-12 of the algorithm, for each variable  $x$ , we calculate the weights of  $x$  taking the values of 0 and 1 respectively according to Equation 1. Subsequently, in line 13, we determine the probability of assigning each variable the value of 1 in the next iteration in accordance with Equation 2. In line 15, we randomly assign values

to each variable based on the sampling probabilities obtained in the previous step, thereby obtaining a reference assignment. However, this reference assignment may not necessarily satisfy the constraints. Consequently, in lines 16-17, we invoke the *ContextSAT* algorithm mentioned in Section II-B to find a valid solution similar to this reference assignment. These steps continues until a set of  $k$  test cases is constructed. Through this process, we can construct a number of test cases that are highly distinct from the current test suite in each iteration, so as to select the test case that can add the maximum number of new tuples.

### C. More Tuples Per Case

In the full coverage phase, *DivSampCA* needs to ensure that the final test suite covers all tuples that do not violate the constraints. Therefore, after the sampling phase, *DivSampCA* continues to add a small number of test cases to achieve full coverage. Notably, *DivSampCA* introduces an innovative technique called *MTPC*. This technique is designed to cover as many tuples as possible with a single test case, thereby reducing the number of test cases needed to cover the remaining uncovered tuples and minimizing the scale of the final PCA. In the subsequent part of this section, we will discuss the SAT conflict analysis mechanism and the implementation of the *MTPC* method.

Building on the concepts introduced in Section II-A, we delve into the advanced aspects of partial assignment. A partial assignment is a map  $\phi : V \rightarrow \{0, 1\}^{|V|}$ , where  $V$  is a finite set of variables and  $|V|$  represents the cardinality of  $V$ . Given a partial assignment, a SAT solver executes the logical reasoning and backtracking search algorithm. The solver will either extend it to a satisfying assignment or claim that the formula is unsatisfiable under the given partial assignment. Leveraging the characteristic that multiple uncovered tuples can coexist simultaneously within the same valid test case, *DivSampCA* traverses the uncovered tuples. It sequentially adds them to a partial assignment and determines whether this partial assignment can be extended into a complete one. If it can be extended, *DivSampCA* will proceed to add the next uncovered tuple; otherwise, it will skip the current tuple and continue to process the next one. This process will keep going until all uncovered tuples have been processed.

The detailed pseudocode of *MTPC* is listed in Algorithm 3. Before running the algorithm, *DivSampCA* first collects all valid uncovered tuples. Given that most tuples have already been covered during the sampling phase, this process is less time-consuming. In line 3 of the algorithm, *MTPC* first shuffles the order of the uncovered tuples to reduce potential biases that may arise from a specific sequence. From line 4 to line 10, *MTPC* traverses the uncovered tuples and processes them with the assistance of the *DPLL* program using the method described above. In line 11, the partial assignment is extended to a complete test case and then integrated into the test suite. Subsequently, the information about uncovered tuples is

---

### Algorithm 3: *MTPC* ( $F$ )

---

**Input:**  $F$ : Boolean formula in CNF;

**Output:**  $\delta$  : a test case for  $F$ ;

```

1  $\alpha \leftarrow \emptyset$ ;
2  $U \leftarrow$  the set of uncovered tuples;
3  $U \leftarrow \text{shuffle}(U)$ ; // Shuffle the order of the
   uncovered tuples
4 foreach  $\tau$  in  $U$  do
5    $\alpha \leftarrow \alpha \cup \{\tau\}$ ;
6   if  $\text{DPLL}(F, \alpha)$  reports “No Solution” then
7      $\alpha \leftarrow \alpha - \{\tau\}$ ;
8     continue;
9   end
10 end
11 if  $\text{DPLL}(F, \alpha)$  returns a solution  $\delta$  then return  $\delta$ ;
12 return  $\emptyset$ ;
```

---

updated to initiate the next round of traversal, and this process continues until all uncovered tuples are fully covered. Throughout this iterative process, *MTPC* constructs test cases that cover multiple uncovered tuples, thereby reducing the size of the final PCA.

**Example 3.** Consider a scenario where we have six distinct uncovered tuples, denoted as  $\tau_1$  through  $\tau_6$ . Due to system constraints, the following combinations are prohibited:  $(\tau_1, \tau_2)$ ,  $(\tau_1, \tau_3)$ ,  $(\tau_1, \tau_5)$ ,  $(\tau_1, \tau_6)$ ,  $(\tau_2, \tau_3)$  and  $(\tau_2, \tau_6)$ .

Applying the *MTPC* method, we first place  $\tau_1$  into a partial assignment. Since  $\tau_1$  conflicts with  $\tau_2$ ,  $\tau_3$ ,  $\tau_5$ , and  $\tau_6$ , we place  $\tau_4$  into the partial assignment and expand it into a complete test case. Similarly, we identify partial assignments for  $\tau_2$  with  $\tau_5$ , and  $\tau_3$  with  $\tau_6$ . Consequently, we cover them with 3 test cases (assuming they can all be extended to complete assignments).

### D. Discussions

1) *Main similarities and differences between DivSampCA and Baital:* Both *DivSampCA* and *Baital* [30] are based on the greedy algorithm framework and adaptively update literal weights to generate diverse test cases. However, *DivSampCA* generates multiple candidate test cases in each iteration and selects the one that maximizes the number of newly covered tuples, while *Baital* generates only one test case per iteration. Most importantly, *Baital* updates the variable assignment probabilities according to the proportion of valid tuples corresponding to each literal. Consequently, it needs to calculate the number of valid tuples for all literals before sampling, which consumes a great amount of time.

2) *Main similarities and differences between DivSampCA and SamplingCA:* Both *DivSampCA* and *SamplingCA* [37] adopt a two-stage algorithm framework and utilize formal reasoning tools to efficiently handle highly configurable software with a large number of configuration options. However, *SamplingCA* updates the variable

assignment probabilities by balancing the number of 0s and 1s for each variable in the test suite. Additionally, during the full coverage phase, *SamplingCA* covers only one uncovered tuple each time, while *DivSampCA* identifies the coexisting uncovered tuples so as to maximize the number of tuples covered by a single test case.

3) *Main similarities and differences between DivSampCA and LS-Sampling-Plus*: *DivSampCA* and *LS-Sampling-Plus* [35] share the common goal of sampling diverse test cases. However, the problems they aim to solve and their approaches differ. *DivSampCA* is designed to address the PCAG problem, aiming to generate a PCA that can cover all valid tuples. In contrast, *LS-Sampling-Plus* targets the  $t$ -wise coverage maximization problem and focuses on creating test suites with high  $t$ -wise coverage. Moreover, not only do they have different mechanisms for updating sampling probabilities, but they also have differences in the test case selection criteria. *DivSampCA* selects test cases that can maximize the number of covered tuples, whereas *LS-Sampling-Plus* chooses test cases based on an approximate scoring function.

#### IV. EXPERIMENTAL DESIGN

To evaluate our method’s performance, we developed a C++ implementation of *DivSampCA*. It is available for public use and accessible online.<sup>1</sup> Before conducting our experiments, we formulated a series of research questions that our experiments are designed to address.

The target of *DivSampCA* is to efficiently generate small-sized PCAs for testing highly configurable systems. This goal leads us to consider two critical aspects of full coverage sampling algorithms: the size of the PCA and the time required for its generation. Consequently, we study the following research questions:

**RQ1** How does *DivSampCA* perform compared with its competitors in PCA size and generation efficiency?

We assessed our two core technologies through ablation experiments to address the following questions:

**RQ2** Is the core technology *TOAS* capable of obtaining test cases that cover more tuples?

**RQ3** Can the core technology *MTPC* cover the same tuples with fewer test cases?

Furthermore, we are also interested in the impact of hyperparameters on the performance of our method.

**RQ4** What is the impact of hyperparameter  $k$  and  $count$  on the performance of *DivSampCA*?

##### A. Benchmarks

In our experiments, we utilized a set of 124 publicly available PCAG benchmark instances.<sup>2</sup> These instances are derived from publicly accessible feature models and represent the practical application requirements across various domains and scenarios. The PCAG instances employed in this study exhibited significant variation in both

size and constraints. Specifically, the number of variables ranged from 94 to 11,254 and the number of constraint clauses ranged from 190 to 62,183, thereby demonstrating the persuasiveness and universality of our experimental results. Given that these public benchmark instances have been thoroughly analyzed in extensive research [36], [38], [44], they possess high generality and representativeness.

##### B. State-of-the-art Competitors

In this study, we conducted experiments by comparing *DivSampCA* with the state-of-the-art algorithms (*i.e.*, *SamplingCA* [37], *CAmpactor* [43] and *LS-Sampling-Plus* [35]) to address the aforementioned questions.

***SamplingCA*** uses sampling techniques to construct a small test suite and adds a few valid test cases to achieve full coverage. Extensive experimental results demonstrate that the performance of the *SamplingCA* outperforms all other PCAG algorithms, including *AutoCCAG* [29], *ACTS* [49], and *HHSa* [50]. The source code for *SamplingCA* is publicly accessible.<sup>3</sup>

***CAmpactor*** is an innovative local search algorithm designed to generate small-sized PCAs. Taking a PCA as input, it removes a specific test case and then searches for another PCA of the resulting size. This process continues until the budget limit is reached. As reported in the literature [43], *CAmpactor* is capable of producing PCAs that are approximately 45% smaller in size compared to those created by current state-of-the-art methods. The implementation of *CAmpactor* is publicly available.<sup>4</sup>

***LS-Sampling-Plus*** iteratively constructs valid test cases based on its novel local search components. According to the experimental results in the literature [35], *LS-Sampling-Plus* can generate test suites with higher coverage than other leading-edge algorithms, such as *NS* [51], *PLEDGE* [52] and *LS-Sampling* [32]. The source code for *LS-Sampling-Plus* can be accessed digitally.<sup>5</sup>

##### C. Experimental Setup

In this work, all experiments were run on a laptop equipped with an Intel(R) Xeon(R) Gold 6226R CPU, with 503GB of RAM and running Ubuntu 23.10. Considering the probabilistic nature of *DivSampCA* and its competitors, we conducted all experiments 10 times to ensure statistical reliability. Furthermore, we selected *MiniSAT* as the implementation of the *DPLL* algorithm. In the experimental setup for *DivSampCA*, both the hyperparameter  $k$  and  $count$  were set to 100, and for the three competing algorithms, we adopted the hyperparameters recommended in the literature [35], [37], [43].

To address RQ1, we evaluated the performance of *DivSampCA* against *SamplingCA* and *CAmpactor* by comparing the average size of the PCAs and the runtime on benchmark instances (with a cutoff time of 3600 CPU

<sup>1</sup><https://anonymous.4open.science/r/DivSampCA-D152>

<sup>2</sup><https://github.com/edbaranov/feature-model-benchmarks>

<sup>3</sup><https://github.com/chuanluocs/SamplingCA>

<sup>4</sup><https://github.com/chuanluocs/CAmpactor>

<sup>5</sup><https://github.com/chuanluocs/LS-Sampling-Plus>

seconds). In addition, following the generation of PCAs by *DivSampCA*, we subsequently employed *LS-Sampling-Plus* to generate a test suite of the same size, enabling a comparison of pairwise coverage and generation time.

To address RQ2 and RQ3, we conducted an ablation study to evaluate the contributions of the two core technologies of *DivSampCA* to the performance improvement. Firstly, *TOAS* strategy was designed to assist in acquiring test cases that cover more tuples during the sampling phase. To assess this, we compared the pairwise tuple coverage achieved by *TOAS* with two other sampling methods (*i.e.*, *Random Sampling* and *SamplingCA*’s *Sampling*) when generating the same number of test cases. Specifically, across 123 instances excluding *embtoolkit* instance, we generated test cases ranging from 30 to 90 in steps of 10. For the *embtoolkit* instance, the number of test cases ranged from 200 to 900 in increments of 100 due to the requirement of approximately 1000 test cases to achieve full coverage. Secondly, we manipulated the specific sets of remaining uncovered tuples by setting the hyperparameter *count* from 30 to 90 with an increment of 10, and ensured that the content of these tuples remained consistent across all experiments. This allowed us to compare the number of test cases required by *DivSampCA* and *SamplingCA* to cover these identical tuples.

To address RQ4, we varied the values of *k* to be 10, 20, 100, 200, and 500, and the values of *count* to be 20, 50, 100, 200, and 500 in the final experiment. The purpose of this manipulation was to explore the impact of the interactions between different settings of *k* and *count* on the performance of *DivSampCA*.

To rigorously assess the statistical significance of performance between *DivSampCA* and its competitors, we conducted the Wilcoxon signed-rank test [53]. Additionally, we calculated the Vargha-Delaney effect size [54] to quantify the magnitude of these differences. If all the p-values from the Wilcoxon signed-rank tests at the 95% confidence level are less than 0.05, and the corresponding Vargha-Delaney effect sizes are all greater than 0.7, we consider the performance of *DivSampCA* to be superior to its competitors, and this superiority is deemed statistically significant. In Table I, the significant results for *DivSampCA* are underlined.

## V. EXPERIMENTAL RESULTS

### A. Comparisons with State of the Art (RQ1)

Table I and Fig. 1 provide a comprehensive comparison of the performance of *DivSampCA* and its competing algorithms. Given that *Campactor* failed to successfully generate PCAs for *embtoolkit* and *uclinux-config* instances, Table I presents the average size and running time of *DivSampCA*, *SamplingCA* and *Campactor* algorithms on the entire instance set excluding these two instances. Additionally, it presents the average coverage rate and time of *LS-Sampling-Plus* when generating an equal number of test cases as *DivSampCA*. In Fig. 1, different line

styles distinguish the algorithms. Furthermore, the red lines represent the sizes of PCAs and correspond to the left axis, while the blue lines denote the CPU time required for generating PCAs and correspond to the right axis.

TABLE I: The performance of *DivSampCA* (denoted as ‘DivSamp’), *SamplingCA* (denoted as ‘Samp’), *Campactor* (denoted as ‘Camp’), and *LS-Sampling-Plus* (denoted as ‘LSP’) over the entire set of instances, excluding the *embtoolkit* and *uclinux-config* instances.

	<i>DivSamp</i>	<i>Samp</i>	<i>Camp</i>	<i>LSP</i>
average size	<u>79.05</u>	105.01	47.07	99.98%
average time (s)	<u>13.51</u>	18.22	226.07	41.66

The experimental results show that in 123 instances out of 124, the PCAs generated by the *DivSampCA* are significantly smaller than those generated by the *SamplingCA*, with an average reduction of approximately 25%. Meanwhile, in 121 out of 124 instances, *DivSampCA* generates PCAs faster than *SamplingCA*, saving an average of 4.71 seconds. In practical testing scenarios, for most feature models, the time required to execute test suites is far longer than that needed to generate them, with each test case taking several hours or even more to execute [55]. Given that the total duration of the testing process is the sum of the time spent on generating and executing test suites, *DivSampCA* not only saves time in generating test suites compared to *SamplingCA* but also reduces the size of the generated test suites by nearly 26 test cases on average. Consequently, it significantly enhances the efficiency of the testing process. In conclusion, the performance of *DivSampCA* is superior to that of *SamplingCA*.

Regarding *Campactor*, although *DivSampCA* generates larger PCAs, on average 1.68 times the size of those generated by *Campactor*, it only takes 0.06 times the time that *Campactor* needs for generation. Moreover, *Campactor*’s reliance on batch generation methods rather than incremental ones results in low efficiency in actual testing scenarios. Testers have to wait for a considerable amount of time for the entire test suite to be completely generated before starting the testing process, rather than receiving newly generated test cases concurrently during the test execution, which consequently reduces the overall testing efficiency. For *LS-Sampling-Plus*, in 122 out of 124 instances, *DivSampCA* is able to successfully generate test suites with 100% coverage and in less time. In contrast, *LS-Sampling-Plus* fails to achieve full coverage and requires a longer time to do so.

**Response to RQ1:** *DivSampCA* can efficiently generate smaller PCAs compared to the current state-of-the-art PCAG algorithms. Specifically, *DivSampCA* produces PCAs that are roughly 25% smaller and 4.71 seconds faster than *SamplingCA* on average. When compared to *Campactor*, although the PCAs generated by *DivSampCA* are 1.68 times larger, the time required is only 0.06 times that of *Campactor*. As for *LS-Sampling-Plus*, *DivSampCA*

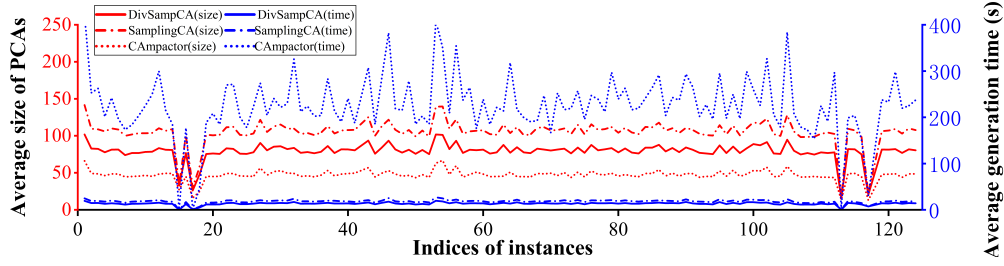


Fig. 1: Performance of *DivSampCA* against *SamplingCA* and *CAmpactor* on the entire instance set.

can more efficiently generate test suites with full coverage in 122 out of 124 instances, whereas *LS-Sampling-Plus* takes more time and fails to achieve full coverage.

### B. Comparisons of Sampling Strategies (RQ2)

Fig. 2 presents the coverage results of *TOAS* and the other two sampling methods, namely *Random Sampling* and *SamplingCA*'s Sampling. The experimental results show that when generating the same number of test cases, *TOAS* consistently achieves a higher average coverage rate over 123 instances compared with the other two methods. This suggests that the *TOAS* technology can effectively sample more diverse sets of test cases, and obtain test cases that cover more tuples.

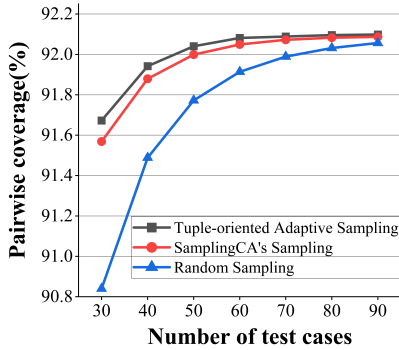


Fig. 2: Average coverage over 123 instances using three sampling methods with an equal number of test cases.

Table II illustrates the proportion of the entire instance set excluding the *embtoolkit* instance in which *TOAS* achieves higher coverage compared to the other two sampling algorithms when generating the same number of test cases. Here, '# tc' represents the number of generated test cases, 'Samp (%)' and 'Rand (%)' denote the proportions in which *TOAS* outperforms *SamplingCA*'s Sampling and *Random Sampling*, respectively. The experimental results indicate that in more than 99% of cases, *TOAS* can sample test cases that cover more tuples than randomly generated test suites. Additionally, compared with *SamplingCA*'s Sampling that focuses on individual literals, *TOAS* technique directly optimizes for tuple coverage, thus enabling it to obtain test cases with higher coverage in at least 95% of the instance set.

TABLE II: Proportion of instances where *TOAS* achieves higher coverage than two other algorithms over the entire instance set except the *embtoolkit* instance.

# tc	30	40	50	60	70	80	90
Samp(%)	96.77	99.19	100	96.77	95.16	96.77	95.97
Rand(%)	100	99.19	99.19	100	99.19	100	100

**Response to RQ2:** The core technology *TOAS* can sample a more diverse set of test cases, thereby enabling the acquisition of test cases that cover more tuples. Compared with the other two methods, *TOAS* consistently achieves a higher average coverage across the entire instance set. Specifically, *TOAS* can generate test suites with higher coverage rates than *Random Sampling* in at least 99% of the instances, and higher coverage rates than *SamplingCA*'s Sampling in at least 95% of the instances.

### C. Effect of the Full Covering Technique (RQ3)

Table III presents the number of test cases required by *DivSampCA* and *SamplingCA*'s full covering method (denoted as 'Samp') to cover the same set of tuples using their respective full covering techniques. The results show that the number of test cases required by *MTPC* is significantly less than that of *SamplingCA*'s full covering method. As the number of tuples increases from approximately 300 to about 1000, the number of test cases required by *SamplingCA*'s full covering method escalates from 3.24 times to 5.52 times that of *MTPC*. This reveals that *MTPC* can indeed identify coexisting tuples among a large number of uncovered tuples, enabling a single test case to cover more tuples simultaneously, thus reducing the number of test cases required.

TABLE III: Number of test cases required to cover the same set of uncovered tuples using *MTPC* and *SamplingCA*'s full covering method, along with their ratio. The number of uncovered tuples is denoted by '# uncov'.

# uncov	321	464	590	721	842	952	1076
Samp	165.74	214.86	253.13	287.89	316.92	343.54	367.63
MTPC	51.13	54.31	56.2	58.95	61.49	64.02	66.56
ratio	3.24	3.96	4.50	4.88	5.15	5.37	5.52

Fig. 3 illustrates the trends in the number of test cases required by two full covering methods as the number of tuples increases. The results indicate that the number of



test cases required by both methods increases with the growth of the number of tuples. However, the growth rate of *MTPC* is significantly slower, merely approximately 1/13 of that of *SamplingCA*'s *full covering* method. This is attributed to the fact that as the number of tuples increases, coexisting tuples become more prevalent. In contrast to *SamplingCA* which directly disregards coexisting tuples, *MTPC* takes them into consideration, enabling it to cover these tuples with fewer test cases. Consequently, appropriately increasing the number of remaining uncovered tuples can enable *MTPC* to play a more pivotal role in reducing the scale of PCAs.

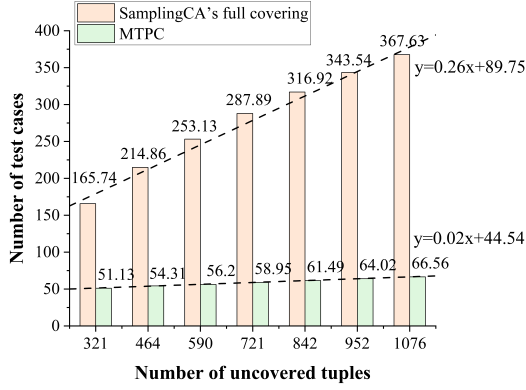
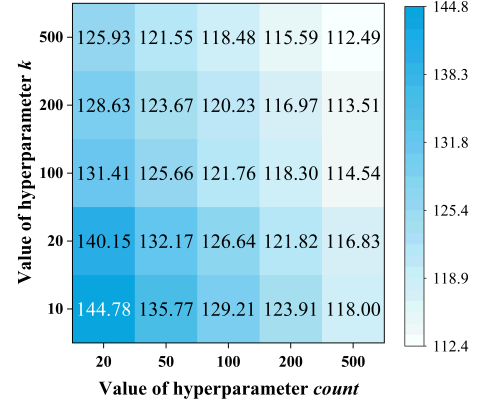


Fig. 3: The number of test cases and their growth rates required by *MTPC* and *SamplingCA*'s *full covering* method to cover the same tuples.

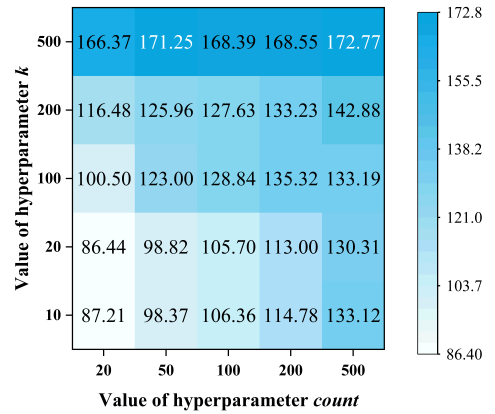
**Response to RQ3:** Compared to *SamplingCA*'s *full covering* method, *MTPC* requires far fewer test cases to cover the same set of uncovered tuples. This trend becomes more pronounced as the number of uncovered tuples increases. In particular, as the number of tuples increases from approximately 300 to 1000, the number of test cases used by *SamplingCA*'s *full covering* method grows from 3.24 times to 5.52 times that of *MTPC*.

#### D. Impacts of Hyper-parameter Settings (RQ4)

Fig. 4 presents the average sizes of PCAs and average times of *DivSampCA* under different hyperparameter settings in the form of a heatmap. In Fig. 4a, the average size varies from 112.49 to 144.78, while in Fig. 4b, the average time ranges from 86.44 to 172.77 seconds. The experimental results show that increasing  $k$  helps to reduce the size of PCAs, but this typically comes at the cost of increased runtime. Moreover, higher values of  $count$  can lead to excessive invocations of the SAT solver, thereby increasing time consumption. On the other hand, lower values of  $count$  may diminish the involvement of *MTPC* during the coverage phase, resulting in larger PCA sizes. In practical testing scenarios, testers can flexibly adjust the values of  $k$  and  $count$  according to specific requirements. The analysis suggests that setting both  $k$  and  $count$  to 100 achieves a trade-off between PCA scale and runtime.



(a)



(b)

Fig. 4: The impacts of hyperparameter on the (a) average size and (b) average time of the *DivSampCA* across the entire instance set. (Darker colors mean worse performance.)

**Response to RQ4:** Increasing the value of  $k$  can effectively enhance the diversity of the test cases, while increasing the value of  $count$  can boost the participation of *MTPC* technique. These parameter settings reduce the size of PCAs but also increase time consumption. Empirically, setting both  $k$  and  $count$  to 100 achieves a balance that optimizes both efficiency and effectiveness.

#### E. Threats to Validity

**Internal Validity.** To reduce the stochasticity of the experimental outcomes, this study conducted multiple experiments. Results show that the coefficient of variation (standard deviation divided by the mean, where a smaller value indicates more stable data) of 114 out of 124 benchmark instances is under 2%, and the highest is 3.1%.

**Construct Validity.** This study adopts the widely recognized performance metrics from prior research. Concretely,

for full coverage algorithms, the size and generation time of PCAs are adopted. For high  $t$ -coverage algorithms,  $t$ -wise coverage is utilized. Moreover, parameter settings for all methods are based on the recommendations in the literature where they were proposed.

**External Validity.** To enhance the generalizability of the research findings, this study employs a diverse and extensive set of representative benchmark instances. These instances cover a broad range of options and constraints, and have been utilized in numerous previous studies.

## VI. RELATED WORK

Combinatorial interaction testing, initially introduced in [56], has been extensively studied in the field of software testing [5]–[7]. In particular, the robust error detection capability of pairwise testing has made it a popular choice for numerous practical applications [6], [14], [57].

The core objective of pairwise testing is to generate a PCA that covers all 2-way interactions. Existing PCAG algorithms are divided into two categories: exact and approximate algorithms. Exact algorithms aim to construct optimal PCAs [8]–[13], while approximate algorithms focus on generating near-optimal PCAs within an acceptable time. The first exact algorithm, *EXACT* [8], was developed by Yan *et al.* Subsequently, Lopez-Escogido *et al.* proposed another exhaustive search method [13] based on *branch-and-bound* techniques. Building on these earlier efforts, a new algorithm named *IPO-MAXSAT* [10] has been proposed, which encodes the problem into the MaxSAT form [58]. Although exact methods can construct optimal PCAs, their excessively high computational complexity restricts their application in practical scenarios.

Approximate algorithms can generally be categorized into three types: greedy, meta-heuristic and sampling algorithms. Greedy algorithms mainly adopt two strategies: *one-test-at-a-time* (OTAT) [14]–[17] and *in-parameter-order* (IPO) [18]–[22]. Bryce *et al.* proposed a well-known OTAT method named *DDA* [17], which can provide a logarithmic guarantee for the size of the PCA. On the other hand, the IPO algorithm proposed by Lei *et al.* [18] starts from a base case and gradually introduces parameters to ensure that the covering array is expanded both horizontally and vertically. This strategy was further developed in the subsequent *IPOG* [19] and *IPOG-F* [20]. Although greedy algorithms can quickly generate PCAs, they may encounter scalability issues [38], making it difficult to produce small PCAs for large-scale problems.

Meta-heuristic algorithms offer diverse approaches for constructing PCAs, including simulated annealing [23]–[26], genetic algorithms [59], [60] and particle swarm optimization [27], [28], [61]. Torres-Jimenez *et al.* proposed a simulated annealing algorithm dubbed *ISA* [25] and enhanced it in their subsequent work [26]. Thereafter, Garvin *et al.* designed another well-known simulated annealing algorithm called *CASA* [62]. In the field of genetic algorithms, Toshiaki *et al.* combined genetic algorithms

with evolutionary strategies to efficiently produce PCAs [59]. Regarding particle swarm optimization approach, Ahmed and Kamal developed a method called *PPSTG* [27]. Based on this strategy, the *PSO-SA* algorithm was proposed to enhance its performance [61]. Meta-heuristic techniques are renowned for their ability to generate small PCAs, but they suffer from issues such as long generation times and inconsistent results.

Sampling methods enable the rapid construction of test cases for systems with numerous configuration options [30]–[37]. Oh *et al.* studied uniform sampling [34], which samples each valid test case with equal probability. Baranov *et al.* proposed an adaptive weighted sampling method in [30], which covers more valid tuples than uniform sampling. Furthermore, Luo *et al.* proposed two recent cutting-edge PCAG algorithms, namely *SamplingCA* [37] and *LS-Sampling* [32] to mitigate the scalability challenge. An enhanced version of *LS-Sampling*, referred to as *LS-Sampling-plus*, was subsequently proposed in [35]. Although sampling methods can quickly generate PCAs, they either fail to ensure full coverage of all tuples or result in overly large PCAs.

In this work, we introduce an innovative tuple-oriented adaptive sampling algorithm named *DivSampCA*. To address the inefficiency of existing PCAG algorithms in generating small-scale PCA problems, *DivSampCA* dynamically adjusts sampling probabilities to generate more diverse test cases and achieves full coverage with fewer test cases. Compared to the state-of-the-art algorithms, *DivSampCA* can effectively produce smaller PCAs. This innovative approach marks an advancement in the realm of sampling-based PCAG algorithms.

## VII. CONCLUSION

In this work, we propose a tuple-oriented adaptive sampling algorithm called *DivSampCA* to address the efficiency issues of existing PCAG algorithms in generating small PCAs. *DivSampCA* incorporates two effective components, *TOAS* and *MTPC*. By sampling diverse test suites based on the proportion of each literal in the uncovered tuples and utilizing individual test cases to cover more tuples, *DivSampCA* efficiently generates smaller PCAs. The results show that considering the interactions between options during the generation of test cases helps guide sampling to avoid the areas already covered by existing test suites, thus facilitating the construction of diverse test cases with less overlap with existing test suites. On the other hand, there are co-existence relationships among different tuples. By leveraging the SAT conflict analysis mechanism, the number of test cases required to cover these tuples can be reduced. However, this work is limited by its greedy addition of test cases that cover the maximum number of tuples. Future work will explore the global benefits of covering the tuples more difficult to cover in the early stages. Additionally, a broader range of test case selection metrics is also an important research direction.

## REFERENCES

- [1] A. Von Rhein, A. Grebhahn, S. Apel, N. Siegmund, D. Beyer, and T. Berger, "Presence-condition simplification in highly configurable systems," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, 2015, pp. 178–188.
- [2] S. Fischer, G. K. Michelon, R. Ramler, L. Linsbauer, and A. Egyed, "Automated test reuse for highly configurable software," *Empirical Software Engineering*, vol. 25, pp. 5295–5332, 2020.
- [3] A. Ali, Y. Hafeez, S. Hussain, and S. Yang, "Role of requirement prioritization technique to improve the quality of highly-configurable systems," *IEEE Access*, vol. 8, pp. 27 549–27 573, 2020.
- [4] C. Nie and H. Leung, "A survey of combinatorial testing," *ACM Comput. Surv.*, vol. 43, no. 2, Feb. 2011.
- [5] I. Dunietz, C. Mallows, A. Iannino, W. Ehrlich, and B. Szablak, "Applying design of experiments to software testing," in *Proceedings of the (19th) International Conference on Software Engineering*, 1997, pp. 205–215.
- [6] D. Kuhn, D. Wallace, and A. Gallo, "Software fault interactions and implications for software testing," *IEEE Transactions on Software Engineering*, vol. 30, no. 6, pp. 418–421, 2004.
- [7] D. Kuhn and M. Reilly, "An investigation of the applicability of design of experiments to software testing," in *27th Annual NASA Goddard/IEEE Software Engineering Workshop, 2002. Proceedings.*, 2002, pp. 91–95.
- [8] J. Yan and J. Zhang, "Backtracking algorithms and search heuristics to generate test suites for combinatorial testing," in *30th Annual International Computer Software and Applications Conference (COMPSAC'06)*, vol. 1, 2006, pp. 385–394.
- [9] J. Bracho-Rios, J. Torres-Jimenez, and E. Rodriguez-Tello, "A new backtracking algorithm for constructing binary covering arrays of variable strength," in *MICAI 2009: Advances in Artificial Intelligence*, A. H. Aguirre, R. M. Borja, and C. A. R. García, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 397–407.
- [10] I. Hiess, L. Kampel, M. Wagner, and D. E. Simos, "IPO-MAXSAT: The in-parameter-order strategy combined with MaxSAT solving for covering array generation," in *2022 24th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC)*, 2022, pp. 71–79.
- [11] B. Hnich, S. D. Prestwich, E. Selensky, and B. M. Smith, "Constraint models for the covering test problem," *Constraints*, vol. 11, no. 2–3, p. 199–219, Jul. 2006.
- [12] D. Lopez-Escogido, J. Torres-Jimenez, E. Rodriguez-Tello, and N. Rangel-Valdez, "Strength two covering arrays construction using a SAT representation," in *MICAI 2008: Advances in Artificial Intelligence*, A. Gelbukh and E. F. Morales, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 44–53.
- [13] J. Torres-Jimenez, I. Izquierdo-Marquez, A. Gonzalez-Gomez, and H. Avila-George, "A branch & bound algorithm to derive a direct construction for binary covering arrays," in *Advances in Artificial Intelligence and Soft Computing*, G. Sidorov and S. N. Galicia-Haro, Eds. Cham: Springer International Publishing, 2015, pp. 158–177.
- [14] D. Cohen, S. Dalal, M. Fredman, and G. Patton, "The AETG system: an approach to testing based on combinatorial design," *IEEE Transactions on Software Engineering*, vol. 23, no. 7, pp. 437–444, 1997.
- [15] S. Abdullah, Z. Soh, and K. Z. Zamli, "Variable-strength interaction for t-way test generation strategy," *International Journal of Advances in Soft Computing & Its Applications*, vol. 5, no. 3, 2013.
- [16] Z. Wang, B. Xu, and C. Nie, "Greedy heuristic algorithms to generate variable strength combinatorial test suite," in *2008 The Eighth International Conference on Quality Software*, 2008, pp. 155–160.
- [17] R. C. Bryce and C. J. Colbourn, "The density algorithm for pairwise interaction testing," *Software Testing, Verification and Reliability*, vol. 17, no. 3, pp. 159–182, 2007.
- [18] Y. Lei and K. Tai, "In-parameter-order: a test generation strategy for pairwise testing," in *Proceedings Third IEEE International High-Assurance Systems Engineering Symposium (Cat. No. 98EX231)*, 1998, pp. 254–261.
- [19] Y. Lei, R. Kacker, D. R. Kuhn, V. Okun, and J. Lawrence, "IPOG: A general strategy for t-way software testing," in *14th Annual IEEE International Conference and Workshops on the Engineering of Computer-Based Systems (ECBS'07)*, 2007, pp. 549–556.
- [20] M. Forbes, J. Lawrence, Y. Lei, R. Kacker, and D. Kuhn, "Refining the in-parameter-order strategy for constructing covering arrays," 2008-09-01 00:09:00 2008.
- [21] K. Kleine and D. E. Simos, "An efficient design and implementation of the in-parameter-order algorithm," *Mathematics in Computer Science*, vol. 12, pp. 51–67, 2018.
- [22] M. Wagner, C. J. Colbourn, and D. E. Simos, "Summary of in-parameter-order strategies for covering perfect hash families," in *2023 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, 2023, pp. 268–270.
- [23] J. Torres-Jimenez, H. Avila-George, and I. Izquierdo-Marquez, "A two-stage algorithm for combinatorial testing," *Optimization Letters*, vol. 11, no. 3, pp. 457–469, 2017.
- [24] J. Torres-Jimenez and E. Rodriguez-Tello, "Simulated annealing for constructing binary covering arrays of variable strength," in *IEEE Congress on Evolutionary Computation*, 2010, pp. 1–8.
- [25] —, "New bounds for binary covering arrays using simulated annealing," *Information Sciences*, vol. 185, no. 1, pp. 137–152, 2012.
- [26] H. Avila-George, J. Torres-Jimenez, and V. Hernández, "Parallel simulated annealing for the covering arrays construction problem," in *Proceedings of the International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA)*. The Steering Committee of The World Congress in Computer Science, Computer ..., 2012, p. 1.
- [27] B. S. Ahmed, K. Z. Zamli, and C. Lim, "The development of a particle swarm based optimization strategy for pairwise testing," *Journal of Artificial Intelligence*, vol. 4, no. 2, pp. 156–165, 2011.
- [28] T. Mahmoud and B. S. Ahmed, "An efficient strategy for covering array construction with fuzzy logic-based adaptive swarm optimization for software testing use," *Expert Systems with Applications*, vol. 42, no. 22, pp. 8753–8765, 2015.
- [29] C. Luo et al., "AutoCCAG: An automated approach to constrained covering array generation," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, 2021, pp. 201–212.
- [30] E. Baranov, A. Legay, and K. S. Meel, "Baital: an adaptive weighted sampling approach for improved t-wise coverage," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2020. New York, NY, USA: Association for Computing Machinery, 2020, p. 1114–1126.
- [31] R. Dutra, K. Laeuffer, J. Bachrach, and K. Sen, "Efficient sampling of SAT solutions for testing," in *Proceedings of the 40th International Conference on Software Engineering*, ser. ICSE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 549–559.
- [32] C. Luo et al., "LS-sampling: an effective local search based sampling approach for achieving high t-wise coverage," in *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2021. New York, NY, USA: Association for Computing Machinery, 2021, p. 1081–1092.
- [33] D.-J. Munoz, J. Oh, M. Pinto, L. Fuentes, and D. Batory, "Uniform random sampling product configurations of feature models that have numerical features," in *Proceedings of the 23rd International Systems and Software Product Line Conference - Volume A*, ser. SPLC '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 289–301.
- [34] J. Oh, P. Gazzillo, and D. Batory, "t-wise coverage by uniform sampling," in *Proceedings of the 23rd International Systems and Software Product Line Conference - Volume A*, ser. SPLC '19.

New York, NY, USA: Association for Computing Machinery, 2019, p. 84–87.

- [35] C. Luo *et al.*, “Solving the t-wise coverage maximum problem via effective and efficient local search-based sampling,” *ACM Trans. Softw. Eng. Methodol.*, vol. 34, no. 1, Dec. 2025.
- [36] —, “Beyond pairwise testing: Advancing 3-wise combinatorial interaction testing for highly configurable systems,” in *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2024. New York, NY, USA: Association for Computing Machinery, 2024, p. 641–653.
- [37] C. Luo, Q. Zhao, S. Cai, H. Zhang, and C. Hu, “SamplingCA: effective and efficient sampling-based pairwise testing for highly configurable software systems,” in *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2022. New York, NY, USA: Association for Computing Machinery, 2022, p. 1185–1197.
- [38] T. Pett *et al.*, “Product sampling for product lines: The scalability challenge,” in *Proceedings of the 23rd International Systems and Software Product Line Conference - Volume A*, ser. SPLC ’19. New York, NY, USA: Association for Computing Machinery, 2019, p. 78–83.
- [39] J. Lin, S. Cai, C. Luo, Q. Lin, and H. Zhang, “Towards more efficient meta-heuristic algorithms for combinatorial test generation,” in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2019. New York, NY, USA: Association for Computing Machinery, 2019, p. 212–222.
- [40] K. Sarkar, C. J. Colbourn, A. de Bonis, and U. Vaccaro, “Partial covering arrays: Algorithms and asymptotics,” in *Combinatorial Algorithms*, V. Mäkinen, S. J. Puglisi, and L. Salmela, Eds. Cham: Springer International Publishing, 2016, pp. 437–448.
- [41] C. Luo *et al.*, “Generating pairwise covering arrays for highly configurable software systems,” in *Proceedings of the 27th ACM International Systems and Software Product Line Conference - Volume A*, ser. SPLC ’23. New York, NY, USA: Association for Computing Machinery, 2023, p. 261–267.
- [42] L. Kampel, M. Leithner, B. Garn, and D. E. Simos, “Problems and algorithms for covering arrays via set covers,” *Theoretical Computer Science*, vol. 800, pp. 90–106, 2019, special issue on Refereed papers from the CAI 2017 conference.
- [43] Q. Zhao *et al.*, “Campactor: A novel and effective local search algorithm for optimizing pairwise covering arrays,” in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 81–93.
- [44] J. H. Liang, V. Ganesh, K. Czarnecki, and V. Raman, “SAT-based analysis of large real-world feature models is easy,” in *Proceedings of the 19th International Conference on Software Product Line*, ser. SPLC ’15. New York, NY, USA: Association for Computing Machinery, 2015, p. 91–100.
- [45] N. Eén and N. Sörensson, “An extensible sat-solver,” in *Theory and Applications of Satisfiability Testing*, E. Giunchiglia and A. Tacchella, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 502–518.
- [46] M. Davis, G. Logemann, and D. Loveland, “A machine program for theorem-proving,” *Commun. ACM*, vol. 5, no. 7, p. 394–397, Jul. 1962.
- [47] N. Manthey, “Coprocessor – a standalone SAT preprocessor,” in *Applications of Declarative Programming and Knowledge Management*, H. Tompits, S. Abreu, J. Oetsch, J. Pührer, D. Seipel, M. Umeda, and A. Wolf, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 297–304.
- [48] G. Seroussi and N. Bshouty, “Vector sets for exhaustive testing of logic circuits,” *IEEE Transactions on Information Theory*, vol. 34, no. 3, pp. 513–522, 1988.
- [49] Y. Jia, M. B. Cohen, M. Harman, and J. Petke, “Learning combinatorial interaction test generation strategies using hyperheuristic search,” in *Proceedings of the 37th International Conference on Software Engineering - Volume 1*, ser. ICSE ’15. IEEE Press, 2015, p. 540–550.
- [50] L. Yu, Y. Lei, M. Nourozborazjany, R. N. Kacker, and D. R. Kuhn, “An efficient algorithm for constraint handling in combinatorial test generation,” in *2013 IEEE Sixth International Conference on Software Testing, Verification and Validation*, 2013, pp. 242–251.
- [51] Y. Xiang, H. Huang, M. Li, S. Li, and X. Yang, “Looking for novelty in search-based software product line testing,” *IEEE Transactions on Software Engineering*, vol. 48, no. 7, pp. 2317–2338, 2022.
- [52] C. Henard *et al.*, “Bypassing the combinatorial explosion: Using similarity to generate and prioritize t-wise test configurations for software product lines,” *IEEE Transactions on Software Engineering*, vol. 40, no. 7, pp. 650–670, 2014.
- [53] W. J. Conover, *Practical nonparametric statistics*. John Wiley & sons, 1999, vol. 350.
- [54] A. Vargha and H. D. Delaney, “A critique and improvement of the cl common language effect size statistics of mcgraw and wong,” *Journal of Educational and Behavioral Statistics*, vol. 25, no. 2, pp. 101–132, 2000.
- [55] Y. Xiang *et al.*, “Automated test suite generation for software product lines based on quality-diversity optimization,” *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 2, Dec. 2023.
- [56] R. Mandl, “Orthogonal latin squares: an application of experiment design to compiler testing,” *Commun. ACM*, vol. 28, no. 10, p. 1054–1058, Oct. 1985.
- [57] A. Hervieu, D. Marijan, A. Gotlieb, and B. Baudry, “Practical minimization of pairwise-covering test configurations using constraint programming,” *Information and Software Technology*, vol. 71, pp. 129–146, 2016.
- [58] C. M. Li and F. Manyá, “MaxSAT, hard and soft constraints,” in *Handbook of satisfiability*. IOS Press, 2021, pp. 903–927.
- [59] T. Shiba, T. Tsuchiya, and T. Kikuno, “Using artificial life techniques to generate test cases for combinatorial testing,” in *Proceedings of the 28th Annual International Computer Software and Applications Conference, 2004. COMPSAC 2004.*, 2004, pp. 72–77 vol.1.
- [60] S. Sabharwal, P. Bansal, and N. Mittal, “Construction of t-way covering arrays using genetic algorithm,” *International Journal of System Assurance Engineering and Management*, vol. 8, no. 2, pp. 264–274, 2017.
- [61] Z. Li, Y. Chen, Y. Song, K. Lu, and J. Shen, “Effective covering array generation using an improved particle swarm optimization,” *IEEE Transactions on Reliability*, vol. 71, no. 1, pp. 284–294, 2022.
- [62] B. J. Garvin, M. B. Cohen, and M. B. Dwyer, “Evaluating improvements to a meta-heuristic search for constrained interaction testing,” *Empirical Software Engineering*, vol. 16, pp. 61–102, 2011.