

# Methods for Evaluating the Similarity of Fuzzers for Model Counters

**Cristian Soare**<sup>1</sup>

# Supervisor: Dr. Anna L. D. Latour<sup>1</sup>

<sup>1</sup>EEMCS, Delft University of Technology, The Netherlands

A Thesis Submitted to EEMCS Faculty Delft University of Technology, In Partial Fulfilment of the Requirements For the Bachelor of Computer Science and Engineering January 27, 2025

Name of the student: Cristian Soare Final project course: CSE3000 Research Project Thesis committee: Dr. Anna L. D. Latour, Dr. Martin Skrodzki

An electronic version of this thesis is available at http://repository.tudelft.nl/.

## Abstract

Model Counting solvers are critical in many domains. One way of validating them is through
fuzzing. However, current fuzzing approaches
lack systematic methods to evaluate how different
test generators compare in bug-triggering behavior.
This paper proposes three methods for evaluating
fuzzer similarity: black-box analysis, differential

8 testing, and white-box analysis.

9 A similarity metric was defined using differential 10 testing, incorporating solver behavior, solve time, 11 and count differences. A case study was conducted 12 on three fuzzers across numerous configurations and four SOTA model counters, followed by an 13 analysis of correlations between CNF feature dis-14 tribution similarities and fuzzer behavioral similar-15 ities. 16

17 The results show moderate correlations between 18 graph structure features (minimum variable node degrees) and fuzzer similarity, while clause bal-19 ance features show negative correlations. However, 20 the method proves inconclusive for selecting dis-21 similar fuzzers due to limited incorrect counts in 22 our dataset, uncertainty about crash causes, and the 23 similarity metric's inherent subjectivity. Further re-24 search is needed, either by expanding the scope of 25 the experiment with more diverse fuzzers, CNF fea-26 tures, and counters, or by pursuing another method 27 recommended by this study. 28

### **29 1** Introduction

Model Counting or #SAT is the counting version of the 30 Boolean Satisfiability (SAT) problem and determines the 31 number of solutions to a Boolean formula in Conjunctive 32 Normal Form (CNF). A solution is defined as a set of as-33 signments of truth values to variables such that the formula 34 evaluates to true. In recent years, many important areas of 35 research have relied more and more on the problem of Model 36 Counting, including areas such as quantitative software veri-37 fication [Teuber and Weigl, 2021, Girol et al., 2021], network 38 reliability [Kabir and Meel, 2023], cryptography [Beck et al., 39 2020], bayesian networks and probabilistic inference [Dar-40 wiche, 2004], and even real-world applications such as reli-41 ability of power grids [Latour et al., 2022]. Numerous such 42 algorithmic problems can be effectively modeled to count the 43 number of solutions to a Boolean formula. Researchers have 44 been more focused on proving that their respective problems 45 can be transformed into #SAT rather than building the neces-46 sary algorithms [Gomes et al., 2021]. Therefore, they require 47 available third-party model counters (#SAT solvers). 48

Using third-party solvers comes with one important issue: 49 they must have a high degree of correctness and efficiency. 50 While it is reasonable to expect model counters to accurately 51 indicate the model count and do this predictably in terms of 52 time, this expectation is not always met. Although new #SAT 53 solvers are built employing testing and validation methods, 54 they also suffer from bugs that are difficult to spot: integer 55 overflow, precision errors, memory issues, and other non-56 deterministic behavior [Brummayer et al., 2010]. For this 57 reason, testing and debugging tools are essential for develop-58 ers building and working with model counters, as they ensure 59 60 reliability and robustness [Brummayer et al., 2010].

Brummayer et al. [2010] argue that **fuzzing** can be used to test SAT solvers and to catch the bugs discussed previously.

This approach can be naturally extended to #SAT solvers as 63 well. Fuzzing works by generating random inputs for the pro-64 gram under test and then running the program with said inputs 65 to check if the program crashes or behaves in an unexpected 66 manner [Zeller et al., 2024]. In the case of #SAT, the output 67 generated by the fuzzers represents a collection of Boolean 68 formulas in CNF form, and an a set of undesired behaviors 69 resulting from running solvers with this input (often charac-70 terized by either a crash or an incorrect model count). 71

Researchers have created tools designed to generate test 72 cases based on the aforementioned principle. However, there 73 is a critical gap in our ability to evaluate them. Specifically, 74 we lack systematic methods to compare their performance 75 across various solvers and to assess the similarities and dif-76 ferences in the test cases they generate. Subsequently, it is 77 difficult for practitioners to choose the most effective fuzzers 78 to determine the correctness of a specific solver. This paper 79 addresses this research gap through the following contribu-80 tions: 81

1. Three distinct methods for evaluating fuzzer similarity in #SAT solvers are introduced, drawing upon methodologies from existing literature. One method is then selected and evaluated through a case study using the SHARPVELVET [Latour and Soos, 2024] project and a collection of instance generators and solvers.

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

123

- 2. Using the selected method:
  - (a) Based on experimental results, CNF features that most effectively characterize fuzzer similarity are explored and analyzed.
  - (b) The method's overall effectiveness in analyzing fuzzer similarity via bug detection patterns is assessed.

The core research question becomes:

What methods can be employed to evaluate and characterize the similarity between different fuzzers for model counting solvers to ensure comprehensive testing coverage?

This overarching question is divided into several specific research questions:

- **RQ1**: What are the key characteristics and limitations 103 of different methodological approaches for evaluating 104 fuzzer similarity in Model Counting solvers? 105
- **RQ2**: Which CNF features most effectively characterize 106 fuzzer similarity when using our approach? 107
- **RQ3**: How effective is our selected method for evaluating fuzzer similarity in Model Counting solvers, and what advantages does it offer compared to alternative approaches?

This paper will proceed with Section 2, which provides 112 more technical definitions for SAT and #SAT. Following that, 113 Section 3 explores related work in the field, offering context 114 and emphasizing the contributions of prior research. Sec-115 tion 4 outlines the proposed methods for evaluating the simi-116 larity of fuzzers, while Section 5 details the methodology and 117 experimental setup for the selected case study. Section 6 ad-118 dresses the transparency, accessibility and scientific integrity 119 part of the study. Section 7 presents the results, limitations, 120 and suggestions for future work. Finally, Section 8 summa-121 rizes the contributions of this research. 122

# 2 Preliminaries

This section contains definitions on Boolean Satisfiability and 124 Model Counting, and important details about SHARPVELVET 125

project. 126

#### 2.1 Boolean Satisfiability 127

Boolean Satisfiability, or SAT, is the problem of determin-128 ing whether a given Boolean formula F (usually in Conjunc-129 tive Normal Form) has at least one satisfying assignment. A 130 formula F in CNF form is defined over a set of Boolean vari-131 ables  $X = \{x_1, x_2, \dots, x_n\}$ , and is represented as a conjunc-132 tion of clauses, where each clause is a disjunction of literals. 133 A literal is either a variable  $x_i$  or its negation  $\neg x_i$ . An assign-134 ment  $\alpha: X \to \{0, 1\}$  is called a satisfying assignment if  $\alpha$ 135 makes F evaluate to True. 136

#### 2.2 Model Counting 137

Model Counting, or #SAT, extends the SAT problem by 138 computing the total number of satisfying assignments of a 139 Boolean formula F. Let Sol(F) denote the set of all satisfy-140 ing assignments of F. The model counting problem seeks to 141 compute |Sol(F)|, the cardinality of Sol(F). 142

In Figure 1, an example problem instance in CNF form and 143 a comparison between the result of the SAT problem and the 144 #SAT problem is shown. Note that the assignment shown for 145 the SAT problem is just one of many correct assignments. 146



Figure 1: An illustration of a boolean formula in CNF form and the difference between SAT and #SAT.

#### 2.3 The SHARPVELVET Project 147

The SHARPVELVET project<sup>1</sup> contains the following separate 148 components: 149

- Instance Generation: The ability to generate CNF in-150 stances in the DIMACS format using multiple genera-151 tors [Latour and Soos, 2024]. 152
- Fuzz Testing: Scripts for executing multiple solvers on 153 the generated instances, reporting discrepancies in their 154 outputs and identifying crashes or incorrect results [La-155 tour and Soos, 2024]. 156
- Verification Tools: Integration with verifiers like CPOG 157 [Bryant et al., 2023] for verifying model counts. 158

SHARPVELVET requires the use of model counters based 159 on the 2024 Model Counting Competition specification<sup>2</sup>. 160 This specification mandates that model counters must accept 161 CNF formulas in DIMACS format and output their results in a 162 163 standardized manner, in order to be later parsed by SharpVel-164 vet and compiled into a report.

#### 3 **Related Work**

This section showcases three important areas relevant to this 166 study: general approaches for fuzzer evaluation, fuzzing tech-167 niques developed for SAT and #SAT solvers, and methods of 168 CNF instance analysis. First, the literature on general fuzzer 169 evaluation is examined, exploring details about experimen-170 tal design. Afterward, fuzzing techniques related to SAT and 171 #SAT solvers are addressed, focusing on black-box fuzzing 172 and the problems of determining the ground truth. The sec-173 tion concludes by discussing methods for characterizing CNF 174 instances, including the established SATZilla feature set and 175 contemporary structural metrics that can be applied to fuzzing 176 analysis. 177

### 3.1 Evaluating Fuzzers

Given that this study focuses on evaluating fuzzers, it is im-179 portant to first look at existing general methods for fuzzer 180 evaluation. Klees et al. [2018] present guidelines for an-181 alyzing new fuzzer innovations, proposing a framework to 182 determine whether a new fuzzer advances the state-of-the-183 art. They emphasize the importance of comparing a pro-184 posed fuzzer against strong baselines, ensuring experiments 185 account for the inherent randomness of fuzzing, and repeat-186 ing tests sufficiently (at least 30 trials) to enable statistical 187 validation [Klees et al., 2018]. More recently, Schloegel et 188 al. [2024] expand on these principles by critically review-189 ing 150 fuzzing papers published between 2018 and 2023. 190 Their analysis revealed widespread deviations from estab-191 lished evaluation guidelines, such as incomplete documen-192 tation of setups and insufficient statistical validation. Both 193 studies highlight the need to report experimental setups, in-194 cluding benchmarks, runtime configurations, and seed selec-195 tion, to support reproducibility and reliable conclusions. 196

Böhme et al. (2022) challenge the reliability of using code 197 coverage to measure fuzzer effectiveness. The study does 198 indicate a strong correlation between code coverage and the 199 number of bugs found. However, only a moderate correlation 200 between coverage and fuzzer ranking is found. This suggests 201 that while higher code coverage may indicate a greater poten-202 tial for uncovering bugs, it does not definitively determine the 203 overall effectiveness of a fuzzer. 204

There are two big differences between the fuzzers de-205 scribed by Klees et al. [2018] and Schloegel et al. [2024] 206 and #SAT fuzzers. Firstly, #SAT fuzzers generate CNF in-207 puts that adhere to a specific format, rather than producing 208 completely random strings. This characteristic complicates 209 the assessment of whether the test set used in the experimen-210 tal evaluation is sufficiently diverse. Furthermore, the em-211 phasis is not on validating whether a particular fuzzer should 212 be classified as state-of-the-art. Instead, comparing the sim-213 ilarity of fuzzers in terms of the bugs they trigger is more 214 relevant. By doing so, it is possible to achieve higher code 215 coverage on solvers by selecting a collection of dissimilar 216 fuzzers. Nonetheless, it remains essential to apply the find-217 ings of Klees et al. [2018] and Schloegel et al. [2024] to en-218 hance the rigor of fuzzer output comparisons, ensuring that 219 the diversity and structural complexity of the generated CNF 220 instances are properly quantified. 221

#### 3.2 Fuzzing Techniques for SAT and #SAT Solvers 222

As mentioned in Section 1, Brummayer et al. [2010] is the 223 cornerstone study that defines the fuzzing SAT and #SAT 224 solvers field. Not only is it the first time fuzzing is proposed 225 as a solution for validating different types of SAT solvers, but 226 it also contributes to one of the earliest implementations of a 227

178

<sup>&</sup>lt;sup>1</sup>https://github.com/meelgroup/SharpVelvet <sup>2</sup>https://mccompetition.org/assets/files/ mccomp\_format\_24.pdf

CNF instance generator. The fuzzing method behind this gen-228 erator is called "black-box fuzzing". The fuzzers themselves 229 cannot interact with the solver-under-test. Instead, they work 230 by only generating syntactically valid but random CNF in-231 puts, which act as inputs to solvers. Given the nature of this 232 method, an important question is: How can we determine if 233 the solver-under-test triggered a bug on a specific instance? 234 This is indeed one important aspect of these fuzzers: they all 235 require a mechanism for deciding on the "ground truth". 236 While crashes are easy to identify, determining whether a 237 model counter has produced an incorrect result presents a 238 more complex challenge. 239

### 240 **3.3** The Structure of Instances

After addressing fuzzing techniques and ways of evaluating 241 fuzzers, another relevant part of this study is analyzing CNF 242 features. Nudelman et al. [2004] are the first to describe the 243 "SATZilla" set of features, which are a collection of features 244 defined to help measure the "hardness" (computational effort) 245 of SAT solvers. This concept is further developed in Xu et 246 al. [2008], who create a solver portfolio which selects from 247 a collection the best counter for each input instance based on 248 an empirical hardness model trained on a set of SATZilla fea-249 tures. Shavit and Hoos [2024] revisited these features and 250 developed a more efficient feature extractor. This updated 251 252 version incorporates modern preprocessing techniques to op-253 timize feature computation for more extensive SAT formulas.

254 The SATZilla set of features has been proven to help predict the computational hardness of SAT solvers. However, 255 an open question remains: how well do these features scale 256 to #SAT solvers? Additionally, do the same features help 257 with fuzzing and bug detection patterns? There are also 258 many other proposed CNF features to consider, such as the 259 treewidth of different graph representations of formulas (pri-260 mal [Dilkas, 2023], consensus, or conflict [Ganian and Szei-261 der, 2021]), structural entropy, as defined by Zhang et al. 262 [2021], and many others. 263

## 264 4 Method Study

This section presents and evaluates three methods for quantifying the similarity of fuzzers for model counters in terms of bug-triggering behavior. The methods were envisioned by documenting standard practices in the fuzzing literature, and merits and weaknesses were analyzed. This section ultimately answers **RQ1**.

#### 271 **4.1 Black-box Analysis**

A black-box oracle evaluates a program solely based on its 272 inputs and outputs, without relying on knowledge of its inter-273 nal workings [Zeller et al., 2024]. In evaluating the genera-274 tors themselves, metrics such as crashes triggered and incor-275 rect results can be used to assess similarities or differences 276 between solvers. By analyzing crash timings, we can deter-277 mine if two fuzzers trigger the same bug, assuming the timing 278 measurements remain consistent and reliable across different 279 executions. For incorrect counts, the analysis becomes more 280 intricate. The number of incorrect model counts per solver 281 or the difference between incorrect and correct ones can help 282 identify which bugs are triggered by which fuzzers. 283

However, this process is not always precise. It often relies on the assumption that similar incorrect counts indicate
the same bug, which may not always hold true [Klees *et al.*,
2018, Böhme *et al.*, 2022]. Improving result metrics can involve gathering extensive samples and applying statistical or
machine-learning methods to identify patterns. For example,

clustering the count differences can reveal bugs with similar behavior. Therefore, this method of evaluating fuzzers only based on solver output requires significant time and effort to produce a comprehensive analysis.

#### 4.2 Differential Testing Approach

Differential testing evaluates program behavior by compar-295 ing solver outputs across multiple implementations [Zeller et 296 al., 2024]. In the context of fuzzer similarity analysis, this 297 approach involves running instances from different fuzzers 298 through multiple solvers and comparing their behavioral pat-299 terns. Like black-box analysis, this method operates solely 300 on solver outputs without requiring internal knowledge of the 301 solvers or fuzzers. 302

The core idea is to characterize the similarity between two 303 fuzzers by analyzing how pairs of their generated instances 304 behave across multiple solvers. A similarity score can be 305 computed for each pair of instances from different fuzzers 306 based on whether they trigger similar solver behaviors (time-307 outs, correct count, incorrect count, crashes, or other out-308 comes). These pairwise similarities can then be aggregated 309 to quantify overall fuzzer similarity. CNF features of the in-310 stances can be analyzed to understand which characteristics 311 contribute to similar behavior patterns. 312

This method is beneficial when traditional black-box anal-313 ysis is impractical - for example, when bugs producing incor-314 rect counts are rare or when precise crash timing analysis is 315 infeasible due to resource sharing or other environmental con-316 straints. However, the approach has significant limitations. 317 The similarity metric is inherently subjective and based on 318 biased assumptions, making it difficult to establish whether 319 two fuzzers truly pinpoint the same solver weaknesses. 320

#### 4.3 White-box Analysis

Unlike analyzing the outputs of different fuzzers in isolation, 322 white-box analysis examines program behavior through sym-323 bolic execution and program code coverage analysis tech-324 niques [Godefroid et al., 2008b]. By exploring execution 325 paths and tracking constraints [Godefroid et al., 2008a], 326 more meaningful insights about test coverage and bug de-327 tection capabilities can be gained compared to black-box ap-328 proaches [Godefroid, 2012]. 329

A method for comparing fuzzers can be implemented by using AFL<sup>3</sup> [Zalewski, 2016] or AFL++ [Fioraldi *et al.*, 331 2020]. The instances produced by existing instance generators like FuzzSAT [Brummayer *et al.*, 2010] can be fed through an AFL-instrumented #SAT solver, allowing execution paths and crashes to be tracked through AFL's shared memory maps [Zeller *et al.*, 2024]. 336

The gathered data can be used to analyze how different 337 fuzzers' outputs engage the solver's functionality, focusing 338 on instances that cause crashes or incorrect counts and identi-339 fying where in the code these issues occur. Using this method, 340 ensuring that two instances trigger the same bug is trivial. 341 Then, CNF features such as SATZilla [Nudelman et al., 2004] 342 can be used to determine which characteristics of the fuzzing 343 input best determine if two generators will trigger the same 344 set of bugs. 345

While this approach has many advantages compared to the previous two, it also has drawbacks. The central problem is that instrumenting a solver causes overhead, which may drastically impact solver performance. This means that the analysis cannot consider solver timeouts or memory exhaustion errors as issues. Another issue is that the instrumented 348

321

<sup>&</sup>lt;sup>3</sup>https://lcamtuf.coredump.cx/afl/

solver setup must first be implemented. Unforeseen bugs re lated to this implementation may arise later in the analysis

and be challenging to identify and address.

# **5 Methodology and Experimental Setup**

Due to resource and time constraints, only the differential 356 testing approach was selected for evaluation. This sec-357 tion focuses on this method and presents a case study using 358 the SHARPVELVET project, a selection of fuzzer configura-359 tions, and state-of-the-art #SAT solvers from the 2024 Model 360 Counting Competition<sup>4</sup>. This method was chosen because 361 it allows an initial assessment of its effectiveness before ex-362 ploring more complex approaches. SHARPVELVET already 363 provides a way of doing differential testing and enables anal-364 ysis of fuzzer behavior patterns even in cases without crashes 365 or incorrect counts, offering insights into similarities between 366 fuzzers beyond traditional bug-triggering methods. 367

A summary of the steps taken to conduct this case study:

- Multiple instance generators were configured with various parameter settings (discussed in Section 5.1 and Section 5.3).
- A set of features to measure the instances generated was
  selected (addressed in Section 5.2).
- 374 3. A collection of diverse state-of-the-art #SAT solvers was 375 chosen (outlined in Section 5.4).
- 4. All instances were fuzzed using instances from the selected generators through the selected #SAT solvers.
- 5. The counter responses for each instance was classified as correct, incorrect, timeout, or crash (further details in Section 5.5).
- 6. Similarity between instances and generators was computed based on the behavioral patterns of the fuzzed solvers (see Section 5.6 and Section 7).
- 7. CNF features and their correlation with the similarity
   computed for each generator was analyzed using statis tical tests (discussed in Section 7).

# 387 5.1 Generator Configurations

This study uses three CNF instance generators: FuzzSAT, 388 PairSAT<sup>5</sup>, and HornSAT<sup>6</sup>, a Horn clause modifier. FuzzSAT, 389 developed by Brummayer et al. [2010] and extended by La-390 tour and Soos [2024], generates formulas by constructing 391 boolean circuits as Direct Acyclic Graphs, randomly select-392 ing operators (AND, OR, XOR, IFF) to connect nodes un-393 til reaching minimum reference counts, then converting to 394 395 CNF via Tseitin transformation. PairSAT, developed by Vuk Jurišić during a parallel research project, creates bipartite 396 graphs representing CNF formulas with controlled clause-to-397 variable ratios and balanced variable distributions. HornSAT, 398 also developed by Vuk Jurišić, takes existing CNF formulas 399 and systematically modifies them to contain varying fraction 400 of Horn clauses (clauses with at most one positive literal). 401

Each base generator (FuzzSAT and PairSAT) is configured
across presumed hardness levels and randomness degrees,
though their actual hardness for model counting remains to
be empirically validated. For each 100% randomness configuration, we select a random instance to create HornSAT
instances. All configurations used and their name can be seen
in Table 1.

Generator	Hardness	Variables	Clause Length	Randomness	
FuzzSAT	Easy	18-30	1-3	0%: Purely structural	
				50%: 45-55% ran-	
				dom clauses	
				100%: 90-100% ran-	
				dom	
	Hard	30-50	4-6	Same as above	
PairSAT	Easy	30-50	1-3	0%: Balanced +	
				fixed (4.0)	
				100%: Unbalanced +	
				variable	
	Hard	55-90	4-6	0%: Balanced +	
				fixed (4.5)	
				100%: Unbalanced +	
				variable	
FuzzSATHORN	Easy	FuzzSAT random easy instance varied to 0-100% Horn clauses FuzzSAT random hard instance, varied to 0-100% Horn clauses			
	Hard				
PairSATHORN	Easy	PairSAT in:	PairSAT instance with same variations PairSAT hard instance with same variations		
	Hard	PairSAT ha			

Table 1: Generator configurations and parameters

The choice of fuzzers and configurations was motivated 409 by the need for a diverse range of instances to evaluate instance similarity. The combination of different generators and 411 randomness parameters enables the analysis of structural patterns across varied formula types, independent of the generators' effectiveness as fuzzers. 414

415

#### 5.2 CNF Features

The feature set used in the analysis is based on the founda-416 tional SATZilla features defined by Nudelman et al, and later 417 reduced by Xu et al. [2008] into a "base" set of features. As 418 mentioned in Section 3, no direct connection has been found 419 between this set of features and the hardness of instances for 420 #SAT. These features were selected because they capture im-421 portant structural patterns within CNF data, which can sub-422 stitute CNF properties, and they are also easy and quick to 423 extract [Shavit and Hoos, 2024]. While other novel CNF fea-424 tures were analyzed and considered, they were excluded from 425 this particular experiment due to the lack of available tools for 426 their extraction in SHARPVELVET. 427

Table 2 describes the feature categories from the original428SATzilla paper [Nudelman *et al.*, 2004].429

Category	Description and Features
Problem Size	Basic formula statistics: variables and clauses before (nvarsOrig, nclausesOrig) and after preprocessing (nvars, nclauses), their ra- tio (vars-clauses-ratio), and preprocessing reductions (reduced- Vars, reducedClauses).
Variable- Clause Graph	Statistics from bipartite graph connecting variables to clauses they appear in. Includes degree distributions for vari- ables (VCG-VAR-mean/min/max/entropy) and clauses (VCG- CLAUSE-mean/min/max/entropy).
Variable Graph	Derived from variable co-occurrence in clauses. Tracks node degree statistics (VG-mean, VG-coeff-variation, VG-min, VG-max).
Clause Graph	Graph where clauses sharing variables are connected. In- cludes degree statistics (CG-mean/min/max/entropy) and clus- tering coefficients measuring local density (cluster-coeff- mean/min/max/entropy).
Balance	Distribution of positive/negative literals in clauses (POSNEG- RATIO-CLAUSE-*) and per variable (POSNEG-RATIO-VAR- *). Also includes clause size proportions (UNARY, BINARY+, TRINARY+).
Horn	Measures Horn formula characteristics through horn-clauses- fraction and variable statistics in Horn clauses (HORNY-VAR- mean/min/max/entropy).

Table 2: Description of SATzilla feature categories used in this analysis.

<sup>&</sup>lt;sup>4</sup>https://zenodo.org/records/14249109

<sup>&</sup>lt;sup>5</sup>https://github.com/Chevuu/PairSAT

<sup>&</sup>lt;sup>6</sup>https://github.com/Chevuu/HornSAT

#### 430 5.3 Instance Reduction Strategy

Initially, 1000 instances per fuzzer configuration were gener-431 ated, resulting in 10,000 total CNF instances (+ 404 Horn-432 SAT instances). The reasoning for using a larger number of 433 instances is that as the quantity of instances generated in-434 creases, the impact of randomness on instance distribution 435 per generator decreases. However, due to time and resource 436 limitations, 100 instances from each generator configuration 437 were used in the fuzzing process. 438

The selection/downsampling process aimed to preserve the 439 distribution of CNF features from the original set. This was 440 achieved through stratified sampling, where the feature space 441 was divided into adaptive bins, and instances were selected 442 proportionally from each bin. Statistical validation confirmed 443 that the reduced dataset maintained the essential characteris-444 tics of the original, as proven by three example feature dis-445 tributions in Figure 2, while reducing computational require-446 ments by 90%. 447



Figure 2: Comparison of three example feature distributions before and after downsampling.

#### 448 5.4 Selected Solvers

The choice of solvers was limited to those accepting DI-449 MACS format instances from the 2024 Model Counting 450 Competition, focusing on traditional, unweighted model 451 counting. Four state-of-the-art solvers were selected based 452 453 on their diverse technical approaches and strong performance. This selection was motivated by the need to ana-454 lyze fuzzing outcomes across different solver implementa-455 tions while keeping the experimental scope computationally 456 tractable. 457

D4 [Lagniez and Marquis, 2017]: A knowledge compilation-based counter using dynamic decomposition with hypergraph partitioning. Achieved strong performance across all sections in the 2024 competition<sup>7</sup>.

- GANAK [Sharma *et al.*, 2019]: Builds on sharpSAT's [Thurley, 2006] component caching architecture with probabilistic component caching and improved heuristics. Currently state-of-the-art in unweighted model counting performance.
- ExactMC [Lai *et al.*, 2021]: Introduces CCDD, a
   Decision-DNNF generalization capturing literal equivalences.
- GPMC [Suzuki *et al.*, 2017]: Extends Glucose 3.0 [Audemard and Simon, 2009] and sharpSAT with preprocessing and tree decomposition-based decision heuristics.

All model counters were run with a timeout of twelve minutes, or 720 seconds. This was chosen because of the limited time and resources available for the study; however, it does balance fuzzing hard instances and the ability to fuzz at least 477 120 instances in 24 hours. 478

#### 5.5 Result Classification

Fuzzing results were characterized by distinct classifications,<br/>as shown in Table 3. The framework distinguishes between<br/>instances with correct counts (including zero for unsatisfiable<br/>instances) and incorrect counts, while also tracking solver-<br/>specific issues like timeouts and crashes.480<br/>481<br/>482

Classification	Details
Correct Count	Solution matches expected count (zero for un- satisfiable instances)
Incorrect Count	Solution differs from expected count
Timeout	Solver exceeded time limit
Crash	Solving process terminated unexpectedly

Table 3: Classification of instances in the testing framework

#### 5.6 Similarity Evaluation

Two fuzzers are evaluated for similarity by analyzing their 486 behavior patterns across multiple solvers. For each in-487 stance generated by a fuzzer, we classify the solver's re-488 sponse into four categories: correct count, incorrect count, 489 timeout, or crash. We then compute a similarity score be-490 tween two fuzzers by comparing these behavioral patterns 491 and their quantitative measures across all their instance pairs 492 and solvers. 493

For two instances  $i_1$ ,  $i_2$  with count  $c_1$ ,  $c_2$  and solve time  $t_1$ ,  $t_2$ , the similarity score  $s(i_1, i_2)$  is calculated as:

	0	if behaviors differ
	$0.5 + 0.5e^{-0.1 t_1 - t_2 }$	if both crashed
$S(i_1, i_2) = \langle$	$0.5 + 0.25e^{-0.01 c_1 - c_2 } +$	if both correct/
	$0.25e^{-0.1 t_1-t_2 }$	incorrect
	P(same behavior)	if both timed out

The exponential decay function provides bounded outputs 494 between 0 and 1 with diminishing rates of change as dif-495 ferences increase. This matches the intuition that the like-496 lihood of similar algorithmic paths decreases exponentially 497 with growing differences in count or solve time. A baseline 498 similarity of 0.5 is assigned for matching behaviors, with the 499 remaining 0.5 weighted equally between count and time dif-500 ferences. The steeper -0.1 decay rate for time reflects how 501 variations of tens of seconds typically indicate distinct algo-502 rithmic paths, while the gentler -0.01 decay for count differ-503 ences accommodates moderate variations that may still rep-504 resent similar underlying behavior. 505

For timeout cases, without having access to the actual count and solve time differences, the exact similarity score cannot be computed. However, rather than assuming that both instances exhibit the same behavior, we compute P(same behavior) as the probability that both instances would exhibit the same behavior (correct or incorrect) based on the sample distributions of their respective generators. 510

The final similarity is calculated between generators as an aggregation of all pairwise similarity scores calculated previously: 514

$$S = \lambda \mu_S + (1 - \lambda)(1 - (s_{max} - s_{min}))$$

where  $\mu_S$  is the mean similarity across all instance pairs and solvers, and  $\lambda = 0.7$  balances between mean similarity and consistency. The second term uses range rather than standard deviation as it directly measures maximum variability 519

485

<sup>&#</sup>x27;https://mccompetition.org/assets/files/ 2024/MC2024\_awards.pdf

in the [0,1]-bounded similarity scores, where higher values

indicate more consistent behavior patterns. This formula-

tion rewards fuzzers that maintain reliable algorithmic pat-terns across all instances rather than just matching on average,

while the range's simple interpretation makes the consistency

525 penalty transparent.

# 526 5.7 Software Environment

The experiments were conducted using an enhanced version of SHARPVELVET [Latour and Soos, 2024], augmented with some additional features<sup>8</sup>. These tools were developed collaboratively by the Research Project Peer Group. The modular architecture preserves the original SHARPVELVET codebase by Latour and Soos [2024] without direct modifications.

# 533 5.8 Hardware Configuration

All experiments were run on the DelftBlue supercomputer [Delft High Performance Computing Centre (DHPC), 2024]. The experiments utilized Phase 1 compute nodes, each equipped with:

• 48-core Intel Xeon processors

- 185 GB RAM
- High-performance interconnect network for parallel computing

# 542 6 Responsible Research

Responsible, accessible and ethical research practices are essential in any scientific study. Therefore, this section discusses such considerations, acknowledging some of the challenges regarding fuzzing model-counting solvers and their impact on the broader research community.

# 548 6.1 Transparency and Reproducibility

An important aspect of conducting experimental research is 549 ensuring that the resulting data is described transparently and 550 that readers of the study can reproduce the experiment in its 551 totality. Given the nature of fuzzing, which employs ran-552 domness throughout the process, including the seed used in 553 the random number generators is essential. Numerous tools 554 555 (generators, counters) were used and can be configured dif-556 ferently. Knowledge of said configuration is required to ensure another person can use these tools similarly. A com-557 panion GitHub repository<sup>8</sup> is therefore available, which in-558 cludes not only any of the aforementioned configurations but 559 also a guide on how to use this information in order to repeat 560 the same experiment. We also include the Jupyter Notebooks 561 used to aggregate the data and generate the figures shown in 562 Section 7. 563

Nevertheless, it is important to acknowledge that the nature 564 of the fuzzing results is non-deterministic and that hardware 565 also plays a role in the possible bugs triggered. For example, 566 running the solvers on slower hardware results in more time-567 outs. Moreover, the operating system and hardware archi-568 tecture dictate arithmetic logic and memory access behavior. 569 Therefore, even with all the information available, a subse-570 quent identical experiment might show different fuzzing be-571

572 havior.

### 6.2 Accessibility

Accessibility is an important consideration, as it directly impacts the study's broader utility and reach. A step taken to achieve this is using colorblind-friendly visualization principles. All figures, tables, and plots are either in black-andwhite or employ a carefully selected colorblind-safe palette called *"Viridis"*. This way, regardless of color impairment, any reader can visualize the figures.

573

601

622

627

A notable limitation to the study's accessibility is its computational requirements. While all experiments were conducted using High-Performance Computing resources [Delft High Performance Computing Centre (DHPC), 2024], these facilities may not be available to researchers aiming to replicate our findings, potentially restricting broader validation of our results. 587

# 6.3 Ethical Considerations and Scientific Integrity 588

Running fuzzers on model-counters has limited ethical implications. The nature of both areas of research is theoretical and does not have direct effects on data privacy or real-world ethical applications. Even if model-counting solvers may be used for unethical purposes, this is derivative and unrelated to the subject of this study: evaluating fuzzers and improving the testing of model counters. 595

Regarding scientific integrity, all software tools used in this paper were utilized under license or with developer permission. They were comprehensively cited, and their licenses are available in their respective repositories and in the main repository for this project.

# 6.4 Generative AI

This research utilized Large Language Models and generative<br/>AI tools to support various aspects of the study. The applica-<br/>tions encompassed:602603604

- Research assistance for brainstorming, source identification, and summarization. For instance, the work by Klees et al. [2018] was identified using ChatGPT's web search functionality. All AI-generated summaries were comprehensively verified against original source materials.
- Help with LaTeX formatting. Example prompt used in Claude 3.5 Sonnet: *"help me make this table smaller in size while keeping the font large. maybe we could also shorten the text in the cells a little to facilitate this"*.
- Software development support via GitHub Copilot, powered by Claude Sonnet 3.5, for implementing shell scripts, integrating feature extraction and verification with generator output in Python, and developing the data aggregation and visualization code.
- Grammarly AI provided writing support through grammar checks and sentence improvements. 620

# 7 Results, Analysis and Future Work

This section presents our experimental results, analyzes the findings, and discusses limitations and suggestions for future research. The experiment was conducted according to the steps described in section 5. 626

# 7.1 Results and Analysis

The analysis was initiated by examining the initial generators and the characteristics of the generated instances. The input data was categorized into the following groups: 630

• Feature Category (Variable Graph Features, Clause 631 Graph Features, etc.) 632

<sup>&</sup>lt;sup>8</sup>https://github.com/Krat-OS/ cse3000-how-to-break-a-solver/tree/ csoare-cse3000-clean

- Base Generator (FuzzSAT, PairSAT, FuzzSATHORN, 633 PairSATHORN) 634
- Presumed Hardness (easy, hard) 635
- Randomness Control (0%, 50%, 100%) 636

While analyzing the distribution of features generated by 637 fuzzers is not the primary focus of this study, it remains im-638 portant to understand the data under evaluation. Figure 3 il-639 640 lustrates the distribution of Clause Graph Features across generators for hard instances with 100% randomness, highlight-641 ing the variations in feature coverage among different gener-642 ation approaches. 643



Figure 3: Distribution of Clause Graph Features for hard instances with 100% randomness.

Due to space constraints and visualization clarity, only rep-644 resentative distributions are presented, with Figure 3 as an il-645 lustrative example of the feature analysis methodology. Key 646 observations from the feature analysis include: 647

- 1. PairSAT and PairSATHORN easy configurations pro-648 duce instances where the feature extractor can reduce all 649 of their original number of clauses and variables to 0. 650 This might be because these instances exhibit duplicate 651 unary clauses. 652
- 2. FuzzSAT shows expected behavior where increased ran-653 domness and hardness correlate with smaller variable-654 to-clause ratios [Nudelman et al., 2004]. 655
- 3. The HornSAT generator maintains consistent Variable 656 Graph and Variable-Clause Graph features with the orig-657 inal input instance while showing more variation in 658 Clause Graph features. 659

660 After fuzzing all of the solvers on the instances described previously, the results in Table 4 were obtained. 661

Counter	Correct	Incorrect	Crash	Timeout
D4	1042	0	102	260
ExactMC	1043	0	0	361
GPMC	1063	0	2	339
Ganak	1086	0	137	181

Table 4: Fuzzing results by model counter.

The incorrect count is determined by comparing each com-662 puted count across all solvers and doing majority voting. As 663 shown in Table 4, the four selected state-of-the-art solvers 664 demonstrated the following behaviors: 665

- All solvers achieved similar correct count rates.
- No incorrect counts were observed across any solver.
- · Crashes and timeouts significantly varied, indicating different performance characteristics across solver implementation. 670

Table 5 presents the fuzzing results categorized by instance 671 generator rather than by solver. 672

Generator	Hardness	Random	Correct	Incorrect	Crash	Timeout
FuzzSAT	easy	0%	397	0	0	3
		50%	400	0	0	0
		100%	400	0	0	0
	hard	0%	157	0	79	164
		50%	400	0	0	0
		100%	400	0	0	0
PairSAT	easy	0%	400	0	0	0
		100%	400	0	0	0
	hard	0%	39	0	55	306
		100%	267	0	71	62
FuzzSAT-	easy	100%	404	0	0	0
HORN	hard	100%	68	0	3	333
PairSAT-	easy	100%	404	0	0	0
HORN	hard	100%	98	0	33	273

Table 5: Fuzzing results by generator.

Already, some interesting patterns regarding instance hard-673 ness emerge. Contrary to our initial assumptions, generators 674 configured for presumably hard instances, which would nor-675 mally result in more timeouts and crashes, did not consis-676 tently produce these results. For example, with FuzzSAT, in-677 creasing the proportion of random clauses in hard instances 678 made them easier to solve and less likely to trigger bugs. This 679 pattern is visible in Figure 4, which shows the average solve 680 time across different generator configurations. Generally, in-681 stances with higher randomness control parameters (either 682 more random clauses in FuzzSAT or unbalanced variable dis-683 tribution with variable clause-to-variable ratios in PairSAT) 684 were solved more quickly. They produced fewer timeouts 685 and crashes than their more structured counterparts. 686



Figure 4: Average solve time by generator.

666 667 668

Also evident from the fuzzing results is that applying the 687 HornSAT modifier to either FuzzSAT or PairSAT hard in-688 stances consistently increases their computational difficulty. 689 The instances produced by each HornSAT modification show 690 higher average solve times and generate more timeouts and 691 crashes than their original counterparts. They also produce 692 instances with a much higher count or number of solutions, 693 as visible in Figure 5. 694



Figure 5: Average count value by generator.

Following the similarity evaluation strategy defined in Section 5.6, every instance from each generator was compared pairwise with all instances from other generators. Figure 6 shows the resulting similarity matrix, defining the similarity score between every two generators.



Figure 6: Pairwise similarity scores between generators, displayed as a triangular matrix to avoid redundant entries.

The matrix highlights distinct yet predictable patterns. In-700 stances generated by presumably easy configurations consis-701 tently achieve high similarity scores (0.7 to 0.9), indicating 702 similar solver behavior and count values. In contrast, in-703 stances generated by harder configurations show lower scores 704 (0.3 to 0.5), reflecting more diverse outcomes, including 705 timeouts and crashes. HornSAT-modified hard instances and 706 PairSAT-hard-0 instances have the lowest similarity scores 707 (< 0.1), which is to be expected because of the longer solve 708 times, higher model counts, and increased timeouts (see Fig-709 ures 4 and 5). 710

In order to answer **RQ2**, we investigated how instance feature distributions predict behavioral similarity between generators. Spearman correlation, which captures monotonic relationships without assuming linearity, was used between feature distribution similarities and behavioral similarity scores. 715 Since our behavioral similarity metric already incorporates solve time and model count differences through specific weights, no additional accounting for these confounding variables was needed. 718



Figure 7: Spearman correlation coefficients between feature distribution similarities and generator behavioral similarity scores. Features are sorted by absolute correlation strength ( $|\rho| > 0.3$ ), showing only statistically significant correlations (p < 0.05).

As shown in Figure 7, only four features showed mod-720 erate absolute Spearman correlations ( $|\rho| > 0.3$ ). Variable 721 graph minimum degrees (VCG-VAR-min, VG-min) showed 722 positive correlations, indicating that generators producing in-723 stances with similar basic graph structures tend to behave 724 similarly. Conversely, minimum positive-to-negative literal 725 ratio (POSNEG-RATIO-CLAUSE-min) and proportion of 726 unary clauses (UNARY) showed negative correlations. 727

These moderate correlations suggest that examining the 728 distribution of these four features could help researchers se-729 lect which fuzzers to run by identifying potentially redun-730 dant behavior without requiring expensive fuzzing campaigns 731 (such as the one conducted in this study). Particularly, fuzzers 732 producing similar VCG-VAR-min and VG-min distributions 733 might be candidates for elimination due to likely behavioral 734 overlap. However, this method proves inconclusive in an-735 swering **RQ3** for two key reasons: first, our dataset could not 736 generate incorrect model counts, weakening the behavioral 737 analysis; second, only a handful of features showed mod-738 erate correlations, and these correlations rely on a subjec-739 tive similarity metric that weighs different behavioral aspects 740 (timeouts, crashes, model counts) based on chosen parame-741 ters rather than empirically derived weights. 742

#### 7.2 Limitations and Future Work

Several limitations of this study suggest directions for future research. Our analysis was constrained by limited computational resources, leading to relatively short timeouts (12 research) and small sample sizes (100 instances per generator). Additionally, the study was limited by the few avail-

743

able model counters and fuzzers, requiring us to use different
 configurations of the same generators to obtain a statistically
 representative dataset.

A significant limitation in our crash analysis is that many reported crashes could have been caused by system-level process termination rather than actual solver bugs. Outof-memory handlers and infrastructure-related issues in the fuzzing setup may have terminated solver processes prematurely, making distinguishing between genuine solver crashes and external termination events difficult.

The absence of incorrect counts in our results indicates the need for more sophisticated instance generators that better explore corner cases. Future work should investigate alternative CNF features beyond the base SATzilla set, such as treewidth measures [Dilkas, 2023] and structural entropy [Zhang *et al.*, 2021].

Our similarity metric, while functionally able to compare 765 fuzzing results, relies on subjective weights for different be-766 havioral aspects, and its black-box nature makes it challeng-767 ing to identify definitively when different errors stem from the 768 same underlying bug. Future work should scale up this exper-769 770 iment with more diverse fuzzers and model counters, develop empirically validated similarity metrics that better capture er-771 ror relationships, and explore the other proposed methodolo-772 gies in Section 4. 773

# 774 8 Conclusion

Although evaluating fuzzers is essential for improving the 775 validation of #SAT solvers, this evaluation is no easy feat. 776 The complex implementation of model counters, combined 777 with limited computational resources and time constraints, 778 makes evaluating fuzzers based on the similarity of bug-779 triggering behavior an essential approach. This study pro-780 posed three methods for achieving this goal while conducting 781 a comprehensive case study using differential testing on three 782 available fuzzers and four state-of-the-art #SAT solvers. 783

The results revealed moderate correlations between cer-784 tain CNF features and fuzzer similarity, particularly in graph 785 structure metrics like variable node degrees and clause polar-786 ity ratios. However, the limited number of features showing 787 significant correlations, combined with the absence of incor-788 rect model counts in our results and uncertainty about crash 789 behaviors, suggests that feature-based similarity alone may 790 791 not be sufficient for selecting a dissimilar subset of fuzzers 792 that maximizes coverage.

While this approach demonstrated promise, the findings indicate the need for more advanced fuzzer evaluation techniques, particularly those incorporating white-box analysis with AFL instrumentation. The proposed methods and results contribute to a growing understanding of effectively evaluating fuzzers for #SAT solvers while working within practical resource constraints.

# References

- [Audemard and Simon, 2009] Gilles Audemard and Laurent Simon. Predicting Learnt Clauses Quality in Modern SAT Solvers. In Proceedings of the 21st International Joint Conference on Artificial Intelligence, IJCAI'09, pages 399–404. Morgan Kaufmann Publishers Inc., 2009.
- [Beck et al., 2020] Gabrielle Beck, Maximilian Zinkus, and Matthew Green. Automating the Development of Chosen Ciphertext Attacks. In 29th USENIX Security Symposium (USENIX Security 20), pages 1821–1837. USENIX Association, 2020.
- [Böhme *et al.*, 2022] Marcel Böhme, László Szekeres, and Jonathan Metzman. On the Reliability of Coverage-Based Fuzzer Benchmarking. In *Proceedings of the 44th International Conference on Software Engineering*, ICSE '22, pages 1621–1633. Association for Computing Machinery, July 2022.
- [Brummayer et al., 2010] Robert Brummayer, Florian Lonsing, and Armin Biere. Automated Testing and Debugging of SAT and QBF Solvers. In Ofer Strichman and Stefan Szeider, editors, *Theory and Applications of Satisfiability Testing – SAT 2010*, pages 44–57. Springer Berlin Heidelberg, 2010.
- [Bryant et al., 2023] Randal E. Bryant, Wojciech Nawrocki, Jeremy Avigad, and Marijn J. H. Heule. Certified Knowledge Compilation with Application to Verified Model Counting. In Meena Mahajan and Friedrich Slivovsky, editors, 26th International Conference on Theory and Applications of Satisfiability Testing (SAT 2023), volume 271 of Leibniz International Proceedings in Informatics (LIPIcs), pages 6:1–6:20. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023.
- [Darwiche, 2004] Adnan Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press, 2004.
- [Delft High Performance Computing Centre (DHPC), 2024] Delft High Performance Computing Centre (DHPC). DelftBlue Supercomputer (Phase 2). https://www.tudelft.nl/dhpc/ark: /44463/DelftBluePhase2, 2024.
- [Dilkas, 2023] Paulius Dilkas. Generating Random Instances of Weighted Model Counting. In Andre A. Cire, editor, Integration of Constraint Programming, Artificial Intelligence, and Operations Research, CPAIOR 2023, pages 395–416. Springer Berlin Heidelberg, 2023.
- [Fioraldi *et al.*, 2020] Andrea Fioraldi, Dominik Maier, Heiko Eißfeldt, and Marc Heuse. {AFL++}: Combining Incremental Steps of Fuzzing Research. In *14th USENIX Workshop on Offensive Technologies*, WOOT 20. USENIX Association, 2020.
- [Ganian and Szeider, 2021] Robert Ganian and Stefan Szeider. New width parameters for SAT and #SAT. *Artificial Intelligence*, 295, June 2021.
- [Girol et al., 2021] Guillaume Girol, Benjamin Farinier, and Sébastien Bardin. Not All Bugs Are Created Equal, But Robust Reachability Can Tell the Difference. In A. Silva and K. R. M. Leino, editors, 33rd International Conference on Computer Aided Verification, CAV 2021, pages 669–693. Cham: Springer International Publishing, 2021.
- [Godefroid *et al.*, 2008a] Patrice Godefroid, Adam Kiezun, and Michael Y. Levin. Grammar-based Whitebox Fuzzing.

In Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '08, pages 206–215. Association for Computing Machinery, 2008.

- [Godefroid *et al.*, 2008b] Patrice Godefroid, Michael Y. Levin, and David Molnar. Automated Whitebox Fuzz Testing. In *Network and Distributed System Security Symposium (NDSS)*, volume 8, pages 151–166. The Internet Society, 2008.
- [Godefroid, 2012] Patrice Godefroid. SAGE: Whitebox Fuzzing for Security Testing. *SAGE*, 10(1), 2012.
- [Gomes et al., 2021] Carla P. Gomes, Ashish Sabharwal, and Bart Selman. Model counting. In Armin Biere, Marijn Heule, and Hans van Maaren, editors, *Handbook of Satisfiability: Second Edition*, Frontiers in Artificial Intelligence and Applications, chapter 25. IOS Press, 2021.
- [Kabir and Meel, 2023] Mohimenul Kabir and Kuldeep S. Meel. A Fast and Accurate ASP Counting Based Network Reliability Estimator. In Ruzica Piskac and Andrei Voronkov, editors, Proceedings of 24th International Conference on Logic for Programming, Artificial Intelligence and Reasoning,, volume 94, pages 270–287, 2023.
- [Klees et al., 2018] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. Evaluating Fuzz Testing. In Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, pages 2123–2138. Association for Computing Machinery, 2018.
- [Lagniez and Marquis, 2017] Jean-Marie Lagniez and Pierre Marquis. An Improved Decision-DNNF Compiler. In Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI'17, pages 667–673, 2017.
- [Lai et al., 2021] Yong Lai, Kuldeep S. Meel, and Roland H. C. Yap. The Power of Literal Equivalence in Model Counting. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35(5), pages 3851–3859, 2021.
- [Latour and Soos, 2024] Anna L.D. Latour and Mate Soos. SharpVelvet, 2024.
- [Latour et al., 2022] Anna L. D. Latour, Behrouz Babaki, Daniël Fokkinga, Marie Anastacio, Holger H. Hoos, and Siegfried Nijssen. Exact Stochastic Constraint Optimization with Applications in Network Analysis. Artificial Intelligence, 304:103650, 2022.
- [Nudelman et al., 2004] Eugene Nudelman, Kevin Leyton-Brown, Holger H. Hoos, Alex Devkar, and Yoav Shoham. Understanding Random SAT: Beyond the Clauses-to-Variables Ratio. In Mark Wallace, editor, *Principles* and Practice of Constraint Programming, CP 2004, pages 438–452. Springer Berlin Heidelberg, 2004.
- [Schloegel et al., 2024] Moritz Schloegel, Nils Bars, Nico Schiller, Lukas Bernhard, Tobias Scharnowski, Addison Crump, Arash Ale-Ebrahim, Nicolai Bissantz, Marius Muench, and Thorsten Holz. SoK: Prudent Evaluation Practices for Fuzzing. In 2024 IEEE Symposium on Security and Privacy (SP), pages 1974–1993, 2024.
- [Sharma et al., 2019] Shubham Sharma, Subhajit Roy, Mate Soos, and Kuldeep S Meel. GANAK: A Scalable Probabilistic Exact Model Counter. In Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI'19, pages 1169–1176, 2019.
- [Shavit and Hoos, 2024] Hadar Shavit and Holger H. Hoos. Revisiting SATZilla Features in 2024. In Supratik

Chakraborty and Jie-Hong Roland Jiang, editors, 27th International Conference on Theory and Applications of Satisfiability Testing (SAT 2024), Leibniz International Proceedings in Informatics (LIPIcs), pages 27:1–27:26. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024.

- [Suzuki et al., 2017] Ryosuke Suzuki, Kenji Hashimoto, and Masahiko Sakai. Improvement of Projected Model Counting Solver with Component Decomposition using SAT Solving in Components. Technical report, JSAI Technical Report, SIG-FPAI-506-07, 2017.
- [Teuber and Weigl, 2021] Samuel Teuber and Alexander Weigl. Quantifying Software Reliability via Model-Counting. In Alessandro Abate and Andrea Marin, editors, 18th International Conference on Quantitative Evaluation of Systems, QEST 2021, pages 59–79. Springer International Publishing, 2021.
- [Thurley, 2006] Marc Thurley. sharpSAT Counting Models with Advanced Component Caching and Implicit BCP. In Armin Biere and Carla P. Gomes, editors, 9th International Conference on Theory and Applications of Satisfiability Testing (SAT 2006), pages 424–429. Springer Berlin Heidelberg, 2006.
- [Xu et al., 2008] Lin Xu, Frank Hutter, Holger H. Hoos, and Kevin Leyton-Brown. SATzilla: Portfolio-Based Algorithm Selection for SAT. *Journal of artificial intelligence* research, 32:565–606, 2008.
- [Zalewski, 2016] Michał Zalewski. American Fuzzy Lop - Whitepaper. https://lcamtuf.coredump.cx/ afl/technical\_details.txt, 2016.
- [Zeller et al., 2024] Andreas Zeller, Rahul Gopinath, Marcel Böhme, Gordon Fraser, and Christian Holler. *The Fuzzing Book.* CISPA Helmholtz Center for Information Security, 2024. Retrieved 2024-07-01 16:50:18+02:00.
- [Zhang *et al.*, 2021] Zaijun Zhang, Daoyun Xu, and Jincheng Zhou. A Structural Entropy Measurement Principle of Propositional Formulas in Conjunctive Normal Form. *Entropy*, 23(3):303, 2021.