



ARCH-COMP18 Category Report: Results on the Falsification Benchmarks

Adel Dokhanchi¹, Shakiba Yaghoubi¹, Bardh Hoxha², Georgios Fainekos¹, Gidon Ernst³, Zhenya Zhang⁴, Paolo Arcaini⁴, Ichiro Hasuo⁴, and Sean Sedwards⁵

¹ School of Computing, Informatics and Decision Systems Engineering,
Arizona State University, Tempe, AZ, U.S.A.

{adokhanc, syaghoub, fainekos}@asu.edu

² Department of Computer Science,
Southern Illinois University, Carbondale, IL, U.S.A.

bhoxha@cs.siu.edu

³ Computing and Information Systems,
University of Melbourne, VIC, Australia

gidon.ernst@unimelb.edu.au

⁴ National Institute of Informatics, Tokyo, Japan

{zhangzy, arcaini, hasuo}@nii.ac.jp

⁵ Intelligent Systems Engineering Lab,
University of Waterloo, ON, Canada

sean.sedwards@uwaterloo.ca

1 Introduction

This report presents the outcomes of the 2018 friendly competition in the ARCH workshop [3] for category of the falsification of temporal logic specifications over Cyber-Physical Systems. The results extend those of the competition of the previous year 2017 by including an additional approach to falsification. Namely, this report includes the 2017 outcomes by the tool S-TALiRO [2] and presents new results by the tool FALSTAR [1]. The benchmarks are available on the ARCH website (cps-vo.org/group/ARCH). In this report, we present results on a powertrain model developed by Toyota Technical Center which contains a complex automatic air-fuel control subsystem [11].

2 Falsification Tools

S-TALiRO [5] is a Matlab toolbox that searches for system behaviors that falsify (do not satisfy) specifications presented in Signal Temporal Logic (STL) [12]. It can analyze arbitrary Simulink models [10] or user-defined black box systems, e.g., autonomous vehicles modeled in a robotics simulator [13]. S-TALiRO performs automated randomized test case generation based on stochastic optimization techniques guided by formal requirements in STL. Among the

advantages of the toolbox is the seamless integration inside the Matlab environment, which is widely used in the industry for model-based development. For a recent overview of the S-TALIRO functionality see [6]. The tool is publicly available on-line at [2] under General Public License (GPL).

FALSTAR is an experimental prototype of a falsification tool that explores the idea to construct falsifying inputs incrementally in time, thereby exploiting potential time-causal dependencies in the problem. It implements several algorithms: time-staging [9], a two layered framework combining Monte-Carlo tree search with stochastic optimization [14], and a probabilistic algorithm that adapts to the difficulty of the problem. The latter algorithm was used for this competition. The code is publicly available under the BSD license [1].

3 Benchmark: Powertrain Control

The Powertrain Control benchmark presented in this report was first introduced in [11]. The benchmark provides a high complexity model of an automatic air-fuel control system. It consists of an air-fuel controller and a mean-value engine model. The closed loop system takes two exogenous inputs: the throttle angle θ_{in} and, the engine speed ω . It has 3 continuous-valued states associated with the controller and 5 continuous-valued states associated with the plant. In addition, there are states which are introduced by the variable delay.

The controller has 4 modes of operation: “Startup”, “Normal”, “Power” and “Fault”. Depending on the operation mode, the system should satisfy different requirements. We used a slightly modified version of the requirements presented in Eq. (27) of the paper by Jin et al. [11]. The following specification needs to be satisfied when the system is in the “Normal” mode:

$$\phi_{PB} = \Box_{(\tau_s, T)}((rise(a) \vee fall(a)) \rightarrow \Box_{(\eta, \zeta)}(|\mu| < \beta))$$

where $a = 40$, $rise(a) \equiv (\theta_{in} \leq 8.8^\circ) \wedge \Diamond_{(0, \epsilon)}(\theta_{in} \geq a)$ for a small enough ϵ , $fall(a)$ is defined similarly, $\tau_s = 11$ is the necessary time for the system to enter the “Normal” mode from the “Startup” mode, $T = 50$ is the total simulation time, $\eta = 1$ is the settling time required after a *rise* or *fall* event happens, $\zeta = 5$ is the end of the current time interval in which the input is kept constant, and, finally, μ is the normalized error signal that indicates the error in the value of the state Air/Flow ratio from a reference value.

The formula states that whenever event *rise* or *fall* happens (the antecedent, which is over the input signal), μ should remain in the specified bound after the settling time η , and before other changes are made to the input (after time ζ). The antecedent of the formula is over the input signals of the system. In this report, the acceptable error bound β is reduced to 0.008 to make falsification feasible. Note that abrupt changes in the value of the input signal are acceptable and necessary here to satisfy the antecedent but, frequent changes in the input are not (less than time ζ). As a matter of fact, increasing the frequency of the changes renders the problem less interesting since falsification becomes easier.

4 Experimental Results

The 2017 experiments with S-TALIRO [2] were conducted on a 64-bit Intel Xeon CPU (2.5GHz) with 64-GB RAM and Windows Server 2012 running MATLAB 2015a. For these experiments, we used the following stochastic optimization methods: Simulated Annealing (SA) [4] and Uniform Random (UR) sampling. We remark that all the experiments were performed with

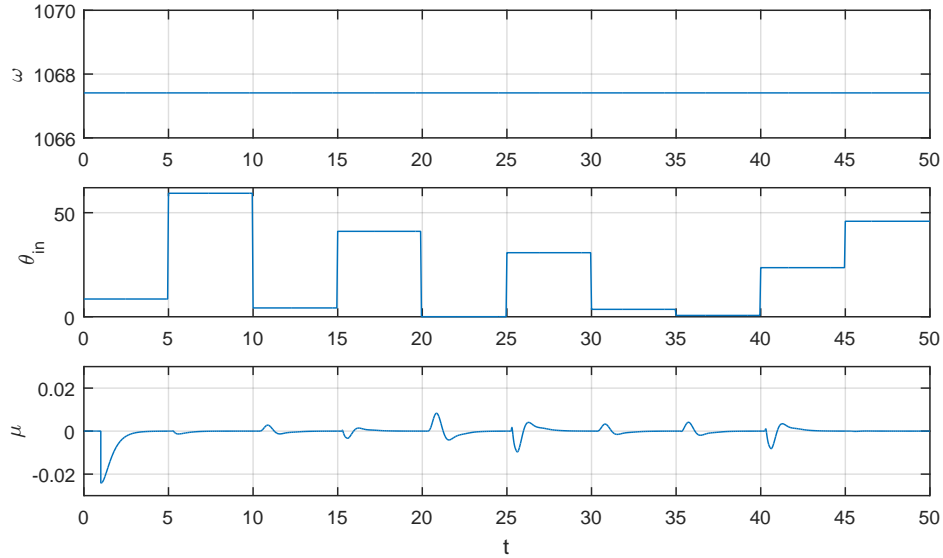


Figure 1: A falsifying piecewise constant input signal (θ_{in}) found by S-TALiRO and the corresponding output trajectory (μ) of the powertrain system for specification ϕ_{PB} . The specification ϕ_{PB} is falsified at time $t = 20.785$, and the robustness value is -6.12×10^{-5} . The input focuses on antecedent falsification up to time $t = 25$. The first 11 sec are ignored based on the requirements in ϕ_{PB} .

the default parameters for each optimization method. It would be expected that further improvements can be achieved by tuning the performance of the optimization algorithms for each benchmark problem. All the benchmark problems are available with the S-TALiRO distribution [2] or from the ARCH workshop repository [3].

The experiments added to the report in 2018 from FALSTAR were conducted on a 64-bit Intel i7-7600U CPU (2.8GHz with 4 cores and 16-GB RAM and Ubuntu 16.04, running MATLAB 2018a). The same Simulink model was used so that the results are comparable. We measure the success rate of finding a falsifying input as a primary indicator of the performance of the tools. Moreover, we state the number of simulations required for success. Note that the computational overhead of the falsification tools is negligible in comparison to running the simulations, so that the number of simulations is indicative of the performance of the tools.

We compare results for the following falsification algorithms.

One is a general falsification algorithm using SA where the optimizer minimizes the robustness value with respect to the given STL specification. This is the standard method used in S-TALiRO. The second one is Vacuity Aware Falsification (VAF) [7]. In VAF, for reactive specifications, as a first step in falsification, we attempt to satisfy the antecedent and, then, falsify the specification. In this report, we review last year's results [8], which demonstrate that S-TALiRO can search over complex constraint input spaces. The S-TALiRO VAF is publicly available [2]. Since the antecedent can be satisfied at any time after τ_s , in our S-TALiRO implementation, in general, we attempt to satisfy the antecedent in a fraction of T ($T/2$ here) so that there is enough time in the future to falsify the whole formula (even though in this particular benchmark this may not be of consequence).

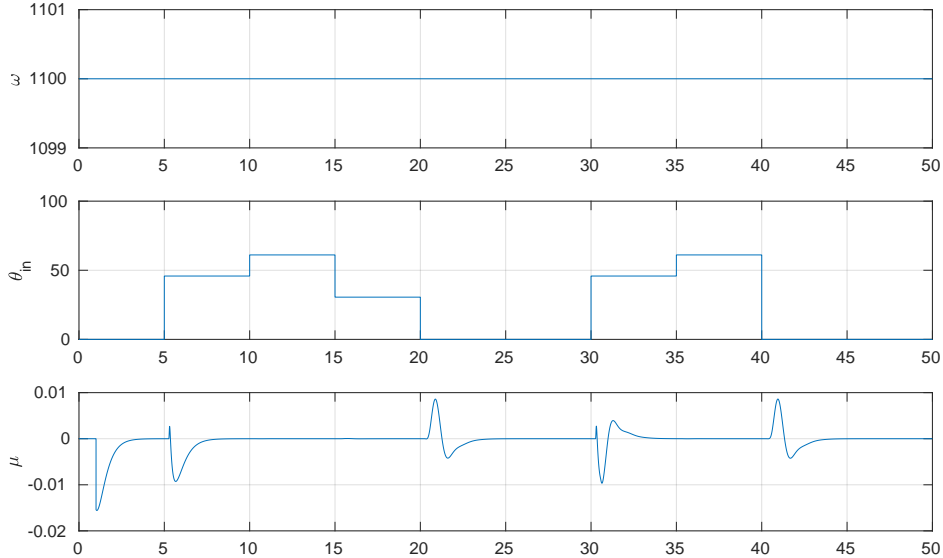


Figure 2: A falsifying trace found by FALSTAR corresponding to the trial with Min Tests = 1 in Table 1. The specification ϕ_{PB} is falsified by the three spikes after $t = 11$.

The third is a novel “Adaptive” algorithm implemented in FALSTAR that constructs input signals incrementally by constructing a tree on the search space. The key feature of the algorithm is that it tries “obvious” input signals first, i.e., such that they take extreme input values $\theta = 0.0$ or $\theta = 61.1$. Only if such values turn out not to be useful, the search gradually switches to more fine-grained choices. Moreover, the algorithm is biased towards extending prefixes that have lead to less robust traces previously.

We used 50 runs (experiments) for each algorithm with 100 tests (simulations) for each run. The experimental results are presented in Table 1. A sample falsifying input and trajectory for S-TALiRO is shown in Fig. 1, and one for FALSTAR is shown in Fig. 2. In the tables, “Min Tests” indicates the minimum number of tests in the case of falsification, while “Min Rob.” indicates the minimum best robustness values achieved for the cases without falsification (not applicable for the adaptive algorithm as it always finds a trace with negative robustness). This gives an idea on how close these cases were to falsification.

Table 1: General Falsification

Optim.	Fals	Min	Max	Avg Tests	Min	Max	Avg Rob.
UR	7/50	18	93	52	1.7×10^{-5}	0.0035	8.81×10^{-4}
SA	9/50	13	83	50	3.54×10^{-5}	0.0042	0.0012
VAF-UR	9/50	12	96	63	3.4×10^{-6}	0.003	0.00086
VAF-SA	29/50	7	95	39	2.38×10^{-6}	0.0043	0.0013
Adaptive	50/50	1	23	6	n/a	n/a	n/a

5 Conclusions

We have presented some preliminary base results for the falsification competition of the ARCH workshop. The results indicate that black box search based test generation methods do not perform much better than random sampling on this challenging benchmark. On the other hand, utilizing some information on the structure of the specification in VAF can help in at least doubling the rate of falsifications.

Furthermore, the adaptive algorithm shows that random sampling under a strong bias towards extreme solutions can in fact lead to high falsification rates. For the powertrain benchmark, this is not too surprising: A large deviation of the air to fuel rate from the reference value is likely linked to extreme changes in the throttle.

As a main outcome of the 2018 ARCH friendly competition, the powertrain model [11] may be considered as solved in the context of falsification tools, and, hence, the participants of the friendly competition will focus on identifying new benchmark problems for the future iterations of the competition.

Acknowledgement FALSTAR is supported by ERATO HASUO Metamathematics for Systems Design Project (No. JPMJER1603), JST; and Grants-in-Aid No. 15KT0012, JSPS. S-TALIRO has been partially supported by NSF awards 1350420, 1446730 and 1361926, and the NSF I/UCRC Center for Embedded Systems.

References

- [1] FalStar : <https://github.com/ERATOMMSD/falstar>.
- [2] S-TaLiRo : <https://sites.google.com/a/asu.edu/s-taliro/>.
- [3] Workshop on Applied Verification for Continuous and Hybrid Systems (ARCH) <http://cps-vo.org/group/ARCH>.
- [4] H. Abbas, G. Fainekos, S. Sankaranarayanan, F. Ivančić, and A. Gupta. Probabilistic temporal logic falsification of cyber-physical systems. *ACM Trans. Embed. Comput. Syst.*, 12(2s):95:1–95:30, May 2013.
- [5] Y. S. R. Annapureddy, C. Liu, G. E. Fainekos, and S. Sankaranarayanan. S-taliro: A tool for temporal logic falsification for hybrid systems. In *Tools and algorithms for the construction and analysis of systems*, volume 6605 of *LNCS*, pages 254–257. Springer, 2011.
- [6] E. Bartocci, J. Deshmukh, A. Donzé, G. Fainekos, O. Maler, D. Nickovic, and S. Sankaranarayanan. Specification-based monitoring of cyber-physical systems: A survey on theory, tools and applications. In *Lectures on Runtime Verification - Introductory and Advanced Topics*, volume 10457 of *LNCS*, pages 128–168. Springer, 2018.
- [7] A. Dokhanchi, S. Yaghoubi, B. Hoxha, and G. Fainekos. Vacuity aware falsification for MTL request-response specifications. In *IEEE International Conference on Automation Science and Engineering*, 2017.
- [8] A. Dokhanchi, S. Yaghoubi, B. Hoxha, and G. E. Fainekos. ARCH-COMP17 category report: Preliminary results on the falsification benchmarks. In *ARCH17. 4th International Workshop on Applied Verification of Continuous and Hybrid Systems, collocated with Cyber-Physical Systems Week (CPSWeek) on April 17, 2017 in Pittsburgh, PA, USA*, pages 170–174, 2017.
- [9] G. Ernst, I. Hasuo, Z. Zhang, and S. Sedwards. Time-staging enhancement of hybrid system falsification. In *Symbolic and Numerical Methods for Reachability Analysis (SNR)*, EPTCS, 2018.
- [10] G. Fainekos, S. Sankaranarayanan, K. Ueda, and H. Yazarel. Verification of automotive control applications using s-taliro. In *Proceedings of the American Control Conference*, 2012.

- [11] X. Jin, J. V. Deshmukh, J. Kapinski, K. Ueda, and K. Butts. Powertrain control verification benchmark. In *Proceedings of the 17th international conference on Hybrid systems: computation and control*, pages 253–262. ACM, 2014.
- [12] O. Maler and D. Nickovic. Monitoring temporal properties of continuous signals. In *Proceedings of FORMATS-FTRTFT*, volume 3253 of *LNCS*, pages 152–166, 2004.
- [13] C. E. Tuncali, G. Fainekos, H. Ito, and J. Kapinski. Simulation-based adversarial test generation for autonomous vehicles with machine learning components. In *IEEE Intelligent Vehicles Symposium (IV)*, 2018.
- [14] Z. Zhang, I. Hasuo, G. Ernst, and S. Sedwards. Two-layered falsification of hybrid systems guided by monte carlo tree search. Preprint, <http://arxiv.org/abs/1803.06276>, 2018.