

Fast Falsification of Hybrid Systems using Probabilistically Adaptive Input

Gidon Ernst

University of Melbourne, Australia
gidon.ernst@unimelb.edu.au

Zhenya Zhang

National Institute of Informatics, Tokyo, Japan
zhangzy@nii.ac.jp

Sean Sedwards

University of Waterloo, Canada
sean.sedwards@uwaterloo.ca

Ichiro Hasuo

National Institute of Informatics, Tokyo, Japan
hasuo@nii.ac.jp

ABSTRACT

We present an algorithm that quickly finds falsifying inputs for hybrid systems, i.e., inputs that steer the system towards violation of a given temporal logic requirement. Our method is based on a probabilistically directed search of an increasingly fine grained spatial and temporal discretization of the input space. A key feature is that it adapts to the difficulty of a problem at hand, specifically to the local complexity of each input segment, as needed for falsification. In experiments with standard benchmarks, our approach consistently outperforms existing techniques by a significant margin. In recognition of the way it works and to distinguish it from previous work, we describe our method as a “Las Vegas tree search”.

CCS CONCEPTS

• **Computer systems organization** → **Embedded and cyber-physical systems**; • **Mathematics of computing** → **Stochastic control and optimization**; • **Theory of computation** → **Random search heuristics**;

KEYWORDS

cyber-physical system, hybrid system, testing, falsification, stochastic optimization, temporal logic, Las Vegas tree search

1 INTRODUCTION

The falsification problem we consider seeks a (time-bounded) input signal that causes a hybrid system model to violate a given temporal logic specification. A popular way to address this is to first construct a “score function” [20] that quantifies how much of the specification has been satisfied during the course of an execution. The falsification can then be treated as an optimization problem, which can be solved using standard algorithms. This approach, especially using “robustness” [16] as the score function, has been successfully applied, resulting in a number of now mature tools [4, 11].

Despite its apparent success, the commonly used robustness semantics [16] is in general not a perfect optimization function. Greedy hill climbing may lead to local optima, hence robustness is only an heuristic score function [21] with respect to the falsification problem. In practice, standard optimization algorithms overcome this limitation by including stochastic exploration. The most sophisticated of these can also model the dynamics of the system (e.g., [2]), in order to estimate the most productive direction of input signal space to explore. There is, however, “no free lunch” [30],

and high performance general purpose optimization algorithms are not necessarily the best choice. For example, such algorithms often optimize with respect to the entire input trace, without exploiting the time causality of the problem, i.e., the fact that a good trace (one that eventually falsifies the property) may be dependent on a good trace prefix.

The contribution of this paper is a fast randomized falsification algorithm (section 3.1) that exploits the time-causal structure of the problem and that adapts to local complexity. In common with alternative approaches, our algorithm searches a discretized space of input signals, but in our case the search space also includes multiple levels of spatial and temporal granularity (section 3.2). The additional complexity is mitigated by an efficient tree search that probabilistically balances exploration and exploitation (section 3.3).

The performance of our algorithm benefits from the heuristic idea to explore simple (coarse granularity) inputs first, then gradually switch to more complex inputs that include finer granularity. Importantly, the finer granularity tends only to be added where it is needed, thus avoiding the exponential penalty of searching the entire input space at the finer granularity. While it is always possible to construct pathological problem instances, we find that our approach is very effective on benchmarks from the literature, including the Powertrain benchmark [22], which is considered a difficult challenge for falsification. Our experimental results (section 4) demonstrate that our algorithm can consistently beat the best existing methods, in terms of speed and reliability of finding a falsifying input.

2 PRELIMINARIES

In this work we represent a deterministic black-box system model as an input/output function $\mathcal{M}: ([0, T] \rightarrow \mathbb{R}^n) \rightarrow ([0, T] \rightarrow \mathbb{R}^m)$. In general, \mathcal{M} comprises continuous dynamics with discontinuities. \mathcal{M} takes a time-bounded, real-valued input signal $\mathbf{u}: [0, T] \rightarrow \mathbb{R}^n$ of length $|\mathbf{u}| = T$ and transforms it to a time bounded output signal $\mathbf{y}: [0, T] \rightarrow \mathbb{R}^m$ of the same length, but potentially different dimensionality. The dimension n of the input indicates that at each moment $t \leq T$ within the time horizon T , the value $\mathbf{u}(t) \in \mathbb{R}^n$ of the input is an n -dimensional real vector (analogously for the output).

We denote by $\mathbf{u}_1\mathbf{u}_2$ the concatenation of signals \mathbf{u}_1 and \mathbf{u}_2 that have the same dimensions, such that $\mathbf{u}_1\mathbf{u}_2: [0, T_1 + T_2] \rightarrow \mathbb{R}^n$. Concatenation of more than two signals follows naturally and is denoted $\mathbf{u}_1\mathbf{u}_2\mathbf{u}_3 \dots$. A constant input signal segment is written (t, v) ,

where t is a time duration and $v \in \mathbb{R}^n$ is a vector of input values. A piecewise constant input signal is the concatenation of such segments.

In this work we adopt the syntax and robustness semantics of STL defined in [13]. The syntax of an STL formula is thus given by

$$\varphi = \neg\varphi \mid \varphi \vee \varphi \mid \varphi \wedge \varphi \mid \varphi \cup_I \varphi \mid \square_I \varphi \mid \diamond_I \varphi \mid \mu, \quad (1)$$

where the logical connectives and temporal operators have their usual Boolean interpretations and equivalences, I is the interval of time over which the temporal operators range, and atomic formulas $\mu \equiv f(x_1, \dots, x_m) \geq 0$ are predicates over the spatial dimensions of a trace. The robustness of trace \mathbf{y} with respect to formula φ , denoted $\rho(\varphi, \mathbf{y})$, is calculated inductively according to the following robustness semantics, using the equivalence $\rho(\varphi, \mathbf{y}) \equiv \rho(\varphi, \mathbf{y}, 0)$.

$$\begin{aligned} \rho(\mu, \mathbf{y}, t) &= f(x_1[t], \dots, x_m[t]), \text{ for } \mu \equiv f(x_1, \dots, x_m) \geq 0 \\ \rho(\neg\varphi, \mathbf{y}, t) &= -\rho(\varphi, \mathbf{y}, t) \\ \rho(\varphi_1 \vee \varphi_2, \mathbf{y}, t) &= \max(\rho(\varphi_1, \mathbf{y}, t), \rho(\varphi_2, \mathbf{y}, t)) \\ \rho(\varphi_1 \wedge \varphi_2, \mathbf{y}, t) &= \min(\rho(\varphi_1, \mathbf{y}, t), \rho(\varphi_2, \mathbf{y}, t)) \\ \rho(\varphi_1 \cup_I \varphi_2, \mathbf{y}, t) &= \max_{t' \in t+I} \left(\min_{t'' \in [t, t']} (\rho(\varphi_1, \mathbf{y}, t'')), \min(\rho(\varphi_2, \mathbf{y}, t')) \right) \\ \rho(\diamond_I \varphi, \mathbf{y}, t) &= \max_{t' \in t+I} (\rho(\varphi, \mathbf{y}, t')) \\ \rho(\square_I \varphi, \mathbf{y}, t) &= \min_{t' \in t+I} (\rho(\varphi, \mathbf{y}, t')) \end{aligned}$$

An important characteristic of the robustness semantics is that it is faithful to standard boolean satisfaction, such that

$$\rho(\varphi, \mathbf{y}) > 0 \implies \mathbf{y} \models \varphi \quad \text{and} \quad \rho(\varphi, \mathbf{y}) < 0 \implies \mathbf{y} \not\models \varphi. \quad (2)$$

Together, these equations justify using the robustness semantics $\rho(\varphi, \mathcal{M}(\mathbf{u}))$ to detect whether an input \mathbf{u} corresponds to the violation of a requirement φ . This correspondence is exploited to find such falsifying inputs through global hill-climbing optimization:

$$\text{Find } \mathbf{u}^* = \arg \min_{\mathbf{u} \in ([0, T] \rightarrow \mathbb{R}^n)} \rho(\varphi, \mathcal{M}(\mathbf{u})) \text{ such that } \rho(\varphi, \mathcal{M}(\mathbf{u}^*)) < 0. \quad (3)$$

Of course, finding an adequate falsifying input \mathbf{u}^* is generally hard and subject to the limitations of the specific optimization algorithm used.

Sound approximations of the lower and upper bounds of the robustness of a prefix \mathbf{y} can sometimes be used to short-cut the search. We thus define lower and upper bounds in the following way.

$$\text{Lower: } \underline{\rho}(\varphi, \mathbf{y}) = \min_{\mathbf{y}'} \rho(\varphi, \mathbf{y}\mathbf{y}') \quad \text{Upper: } \bar{\rho}(\varphi, \mathbf{y}) = \max_{\mathbf{y}'} \rho(\varphi, \mathbf{y}\mathbf{y}') \quad (4)$$

A lower bound $\underline{\rho}(\varphi, \mathcal{M}(\mathbf{u})) > 0$ can be used to detect that a prefix cannot be extended to a falsifying trace (e.g., after the deadline for a harmful event has passed). An upper bound $\bar{\rho}(\varphi, \mathcal{M}(\mathbf{u})) < 0$ similarly implies $\mathcal{M}(\mathbf{u}\mathbf{u}') \not\models \varphi$ for all \mathbf{u}' , concluding that input \mathbf{u} is already a witness for falsification (e.g., a limit is already exceeded).

3 APPROACH

We wish to solve the following falsification problem efficiently:

$$\text{Find } \mathbf{u}^* \text{ such that } \bar{\rho}(\varphi, \mathcal{M}(\mathbf{u}^*)) < 0. \quad (5)$$

Our approach is to repeatedly construct input signals $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_3 \dots$, where \mathbf{u}_i is drawn from a predetermined search space of candidate

input segments, \mathcal{A} . The choice is probabilistic, according to a distribution \mathcal{D} that determines the search strategy, i.e., which inputs are likely to be tried next given a partially explored search space. The construction of each input is done incrementally, to take advantage of the potential short cuts described at the end of section 2.

Algorithm 1 “adaptive Las Vegas Tree Search” (aLVTS) codifies the high level functionality of this probabilistic approach, described in detail in section 3.1.

The effectiveness of our algorithm in practice comes from the particular choices of \mathcal{A} and \mathcal{D} , which let the search gradually *adapt* to the difficulty of the problem. The set \mathcal{A} (defined in section 3.2) contains input segments of diverse granularity, which intuitively corresponds to how precise the input must be in order to find a falsifying trace. The distribution \mathcal{D} (defined in section 3.3) initially assigns high probabilities to the “coarsest” input segments in \mathcal{A} . Coarse here means that the segments tend to be long in relation to the time horizon T and large in relation to the extrema of the input space. The algorithm probabilistically balances exploration and exploitation of segments, but as the coarser segments become fully explored, and the property has not been falsified, the algorithm gradually switches to finer-grained segments.

3.1 Algorithm

Algorithm 1: Adaptive Las Vegas Tree Search (aLVTS)

Input:
system model $\mathcal{M} : \mathbf{u} \rightarrow \mathbf{y}$,
with $\mathbf{u} : [0, t] \rightarrow \mathbb{R}^n$ and $\mathbf{y} : [0, t] \rightarrow \mathbb{R}^m$
time-bounded specification ϕ
set of all possible input trace segments \mathcal{A}

Output:
 \mathbf{u} such that $\mathcal{M}(\mathbf{u}\mathbf{u}') \not\models \phi$ for all \mathbf{u}' , or
 \perp after timeout or maximum number of iterations

- 1 $unexplored(\mathbf{u}) \leftarrow \mathcal{A}$ for all \mathbf{u}
- 2 $explored(\mathbf{u}) \leftarrow \emptyset$ for all \mathbf{u}
- 3 **repeat**
- 4 $\mathbf{u} \leftarrow \langle \rangle$
- 5 **while** $unexplored(\mathbf{u}) \neq \emptyset \vee explored(\mathbf{u}) \neq \emptyset$ **do**
- 6 sample $\mathbf{u}' \sim \mathcal{D}(\mathbf{u})$
- 7 **if** $\mathbf{u}' \in unexplored(\mathbf{u})$ **then**
- 8 $unexplored(\mathbf{u}) \leftarrow unexplored(\mathbf{u}) \setminus \{\mathbf{u}'\}$
- 9 $\mathbf{y} \leftarrow \mathcal{M}(\mathbf{u}\mathbf{u}')$
- 10 **if** $\bar{\rho}(\phi, \mathbf{y}) < 0$ **then**
- 11 **return** $\mathbf{u}\mathbf{u}'$
- 12 **if** $\underline{\rho}(\phi, \mathbf{y}) > 0$ **then**
- 13 **continue line 3**
- 14 $explored(\mathbf{u}) \leftarrow explored(\mathbf{u}) \cup \{\mathbf{u}'\}$
- 15 $\mathbf{u} \leftarrow \mathbf{u}\mathbf{u}'$
- 16 **until** timeout or maximum number of iterations;
- 17 **return** \perp

Algorithm 1 searches the space of input signals constructed from piecewise constant (over time) segments, which are chosen

at random according to the distribution defined by \mathcal{D} in line 6. This distribution is a function of the numbers of unexplored and explored edges at different levels of granularity, and thus defines the probabilities of exploration, exploitation and adaptation. The precise calculation made by \mathcal{D} is described in section 3.3.

As the search proceeds, the algorithm constructs a tree whose nodes each correspond to a unique input signal prefix. The edges of the tree correspond to the constant segments that make up the input signal. The root node corresponds to time 0 and the empty input signal (line 4).

To each node identified by an input signal prefix \mathbf{u} is associated a set of unexplored edges, $unexplored(\mathbf{u}) \subseteq \mathcal{A}$, that correspond to unexplored input signal segments, and a set of explored edges, $explored(\mathbf{u}) \subseteq \mathcal{A}$, that remain inconclusive with respect to falsification. Initially, all edges are unexplored (line 1 and line 2). Once an edge has been chosen (line 6), the unique signal segment associated to the edge may be appended to the signal prefix associated to the node, to form an extended input signal. If the chosen edge is unexplored, it is removed from the set of unexplored edges (line 8) and the extended input signal \mathbf{uu}' is transformed by the system into an extended output signal (line 9). If the requirement is falsified by the output signal (\mathbf{y} in line 10), the algorithm immediately quits and returns the falsifying input signal (line 11). If the requirement is satisfied, with no possibility of it being falsified by further extensions of the signal (line 12), the algorithm quits the current signal (line 13) and starts a new signal from the root node (line 4). This is the case, in particular, when the length of the signal exceeds the time horizon of the formula as a consequence of the definition of $\underline{\rho}$ in eq. (4). If the requirement is neither falsified nor satisfied, the edge is added to the node's set of explored edges (line 14). Regardless of whether the chosen edge was previously explored or unexplored, if the signal remains inconclusive, the extended input signal becomes the focus of the next iterative step (line 15).

3.2 Definition of \mathcal{A}

Given an input signal segment \mathbf{u}' of length t time units and value $(v_1, \dots, v_n) \in \mathbb{R}^n$, let \underline{v}_i and \bar{v}_i denote the minimum and maximum possible values, respectively, of dimension $i \in \{1, \dots, n\}$. For each integer level $l \in \{0, \dots, l_{\max}\}$, we define the set of possible proportions of the interval $[\underline{v}_i, \bar{v}_i]$ as

$$\mathbf{p}_l = \{(2j+1)/2^l \mid j \in \mathbb{N}_0 \leq (2^l - 1)/2\}.$$

The numerators of all elements are coprime with the denominator, 2^l , hence $\mathbf{p}_i \cap \mathbf{p}_j = \emptyset$ for all $i \neq j$. By definition, \mathbf{p}_0 also includes 0. Hence, $\mathbf{p}_0 = \{0, 1\}$, $\mathbf{p}_1 = \{\frac{1}{2}\}$ and $\mathbf{p}_2 = \{\frac{1}{4}, \frac{3}{4}\}$, etc. The set of possible values of dimension i at level l is thus given by

$$\mathbf{v}_{i,l} = \underline{v}_i + \mathbf{p}_l \times (\bar{v}_i - \underline{v}_i).$$

Rather than making the granularity of each dimension independent, we interpret the value of l as a granularity “budget” that must be distributed among the (non-temporal) dimensions of the input signal. The set of possible per-dimension budget allocations for level l is given by

$$\mathbf{b}_l = \{(b_1, \dots, b_n) \in \mathbb{N}_0^n \mid b_1 + \dots + b_n = l\}.$$

For example, with $n = 2$, $\mathbf{b}_3 = \{(0, 3), (1, 2), (2, 1), (3, 0)\}$. If we denote the set of possible time durations at level l by \mathbf{t}_l , then the

Table 1: Size of \mathcal{A}_l for input dimensionality n and level l given that the size of \mathbf{t}_l is one.

$n \mid l =$	0	1	2	3	4	5	6	7	8	9	10	
2		4	4	9	20	44	96	208	448	960	2048	4352
3		8	12	30	73	174	408	944	2160	4896	11008	24576

set of possible input segments at level l is given by

$$\mathcal{A}_l = \bigcup_{(b_1, \dots, b_n) \in \mathbf{b}_l} \mathbf{t}_l \times \mathbf{v}_{1, b_1} \times \dots \times \mathbf{v}_{n, b_n}.$$

Note that while \mathbf{t}_l is not required here to share the granularity budget, this remains a possibility. Our implementation (section 4.1) actually specifies \mathbf{t}_l by defining a fixed number of control points per level, $(k_0, \dots, k_{l_{\max}})$, such that the $\mathbf{t}_l = \{T/k_l\}$ are singleton sets. The sizes of various \mathcal{A}_l for different choices of n and l , assuming $|\mathbf{t}_l| = 1$, are given in table 1.

In summary, an input segment $\mathbf{u}' = (t, v_1, \dots, v_n) \in \mathcal{A}_l$ has $t \in \mathbf{t}_l$ and corresponding budget allocation $b_1 + \dots + b_n = l$, with the value v_i for each dimension given by $v_i = \underline{v}_i + p_i(\bar{v}_i - \underline{v}_i)$, where $p_i = (2j_i + 1)/2^{b_i}$, for some j_i , defines the proportion between minimum \underline{v}_i and maximum \bar{v}_i .

By construction, $\mathcal{A}_i \cap \mathcal{A}_j = \emptyset$, for all $i \neq j$. Hence, we define

$$\begin{aligned} unexplored_l(\mathbf{u}) &= unexplored(\mathbf{u}) \cap \mathcal{A}_l \text{ and} \\ explored_l(\mathbf{u}) &= explored(\mathbf{u}) \cap \mathcal{A}_l. \end{aligned}$$

The set of all possible input signal segments is given by $\mathcal{A} = \mathcal{A}_0 \cup \mathcal{A}_1 \cup \dots \cup \mathcal{A}_{l_{\max}}$.

Figure 1 depicts the construction of \mathcal{A} for two dimensions. The majority of candidate input points is concentrated on the outer contour, corresponding to an extreme choice for one dimension and a fine-grained choice for the other dimension. While this bias appears extreme, as layers are exhausted, finer-grained choices become more likely. For example, after two points from \mathcal{A}_0 have been tried, all remaining points from both levels in the second panel would be equally probable.

3.3 Definition of \mathcal{D}

The distribution $\mathcal{D}(\mathbf{u})$ is constructed implicitly. First, a granularity level $l \in \{0, \dots, l_{\max}\}$ is chosen at random, with probability in proportion to the fraction of edges remaining at the level, multiplied by an exponentially decreasing scaling factor. Defining the overall weight of level l as

$$w_l = \frac{|unexplored_l(\mathbf{u})| + |explored_l(\mathbf{u})|}{2^l \cdot |\mathcal{A}_l|},$$

level l is chosen with probability $w_l / \sum_{i=0}^{l_{\max}} w_i$.

Having chosen l , one of the following strategies is chosen uniformly at random:

- (1) select $\mathbf{u}' \in unexplored_l(\mathbf{u})$, uniformly at random;
- (2) select $\mathbf{u}' \in explored_l(\mathbf{u})$, uniformly at random;
- (3) select $\mathbf{u}' \in explored_l(\mathbf{u})$, uniformly at random from those that minimise $\rho(\varphi, \mathcal{M}(\mathbf{uu}'))$;

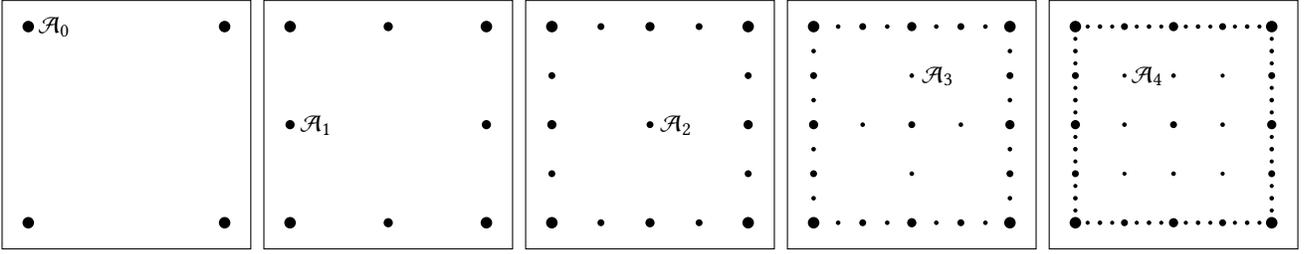


Figure 1: Construction of \mathcal{A} for $n = 2$ and $l_{\max} = 4$, with \mathcal{A}_0 representing the extremes of the two spatial dimensions. Showing $\mathcal{A}_0, \mathcal{A}_0 \cup \mathcal{A}_1, \dots, \mathcal{A}_0 \cup \mathcal{A}_1 \cup \mathcal{A}_2 \cup \mathcal{A}_3 \cup \mathcal{A}_4$. Larger points correspond to more likely values.

- (4) select $\mathbf{u}' \in \text{explored}_l(\mathbf{u})$, uniformly at random from those that minimise $\rho(\varphi, \mathcal{M}(\mathbf{u}\mathbf{u}'\mathbf{u}^*))$, where \mathbf{u}^* denotes any, arbitrary length input signal suffix that has already been explored from $\mathbf{u}\mathbf{u}'$.

Strategy 1 can be considered pure exploration, while strategies 3-4 are three different sorts of exploitation. In the case that $\text{unexplored}_l(\mathbf{u})$ or $\text{explored}_l(\mathbf{u})$ are empty, their corresponding strategies are infeasible and a strategy is chosen uniformly at random from the feasible strategies. If for all $\mathbf{u}' \in \text{explored}(\mathbf{u})$, $\text{explored}(\mathbf{u}\mathbf{u}') = \emptyset$, then strategy 4 is equivalent to strategy 3, but it is *not* infeasible.

3.4 Design Choices

In contrast to Monte Carlo tree search (MCTS), which typically makes a good decision by reaching a goal state multiple times during ployout, our algorithm expects to reach a goal state only once. In place of a statistic about reaching goal states, we therefore adopt a heuristic, namely the lowest robustness. Using this heuristic is typically much more effective than choosing entirely at random, but simply substituting the heuristic value for the statistic in the UCT formula (upper confidence bound applied to trees [25], successfully used with MCTS) is plausible, but nevertheless questionable in our application. Moreover, the deterministic way that UCT balances exploration and exploitation is often not optimal in the context of falsification.

Our design choices are guided by the observation that simple, coarse-grained input signals either immediately solve typical falsification problems, or need only small modifications to do so. Since such signals may be quickly explored exhaustively, our level scaling factors are designed to make this happen with high probability.

Unexplored edges have unknown potential, while every explored edge does not yet lead to falsification, but does not exclude it. The only exception to this occurs when all the edges of a node have been explored and discarded. In this unlikely case the edge leading to such a node exists and can be traversed, but the node and its trace are immediately rejected (line 5).

While not explicit in the presentation of algorithm 1, our approach is deliberately incremental in the evaluation of the system model. In particular, we can re-use partial simulations to take advantage of the fact that traces share common prefixes. Hence, for example, one can associate to every visited \mathbf{u} the terminal state of the simulation that reached it, using this state to initialize a new simulation when subsequently exploring $\mathbf{u}\mathbf{u}'$. This idea also works

for the calculation of robustness. We note, however, that incremental simulations may be impractical. For example, suspending and re-starting Simulink can be more expensive than performing an entire simulation from the start. Under such circumstances, the inner loop of our algorithm can be replaced by a recursion that executes the model on the complete input, returning the trace from the innermost recursion. The checks in lines 11 and 13 can be done on the appropriate prefix of the output returned from the recursive call (which takes place in line 9).

4 EVALUATION

We briefly describe our implementation and the benchmarks, then present the results of an experimental comparison with random sampling and Breach.

4.1 Implementation

We have implemented algorithm 1 in the prototype tool FALSTAR, using the Scala programming language and interfacing to MATLAB/Simulink through a Java API. The code is publicly available on github,¹ including all files and instructions required to re-run our experiments.

The main data structure is an explicit representation of the search tree. Its nodes are labelled by inputs and the scores relevant for strategies 3 and 4 of section 3.3. Similarly, we keep track of the sets unexplored_l and explored_l for each level l .

We do not re-run the simulation for every new input segment, as implied by line 9 of algorithm 1. Instead, we cache the output traces within the tree data structure, alongside their inputs. As mentioned in section 3.4, the cost of pausing and resuming a simulation in Simulink incurs a prohibitively large overhead, hence we have implemented a recursive algorithm that executes exactly one simulation up to the full time horizon T per trial/iteration of the outer loop.

The implementation supports initial parameters with respect to the input signals, such as engine speed ω that should remain constant, by augmenting the first level of the search tree with an additional dimension for each of the parameters.

Currently, the possible time durations t_l for each level l can be specified statically as singleton sets, using a sequence of numbers of control points $(k_0, \dots, k_{l_{\max}})$, such that $t_l = \{T/k_l\}$. The particular choices depend on the requirements of the case study, as noted

¹<https://github.com/ERATOMMSD/falstar>

below. We make use of both progressive subdivisions (increasing k_l with l) and constant subdivision (keeping k_l the same for all l).

We use the algorithm from [12] to compute the robustness of STL formulas efficiently, and the approach from [14] to compute the respective upper and lower bounds.

FALSTAR exposes its functionality through a simple specification and scripting language. In the spirit of the SMT-LIB file format [6], the language is S-expression-based, to facilitate integration with other tools. Using this interface one can declare models, input signal specifications, as well as STL requirements. In addition to solving falsification problems and generating reports, FALSTAR can be used to simulate systems for a given input or to compute robustness values of given STL requirements for specified output traces. To aid comparisons, FALSTAR can also call Breach as an alternative solver, or generate a standalone MATLAB script to run Breach.

4.2 Benchmarks

The benchmarks here are based on automotive Simulink models that are commonly used in the falsification literature. We focus on those benchmarks that have at least one time-varying input, in contrast to just searching for an initial configuration.

Automatic Transmission. This benchmark was proposed in [18] and consists of a Simulink model and its related requirements. The model has two inputs: *throttle* and *brake*. Depending on the speed and engine load, an appropriate gear is automatically selected. Besides the current gear g , the model outputs the speed v and engine ω . We consider the following requirements:

$$\text{AT1} = \square_{[0,30]} v < 120$$

$$\text{AT2}(i) = \square_{[0,30]} (g = i \implies v > 10 \cdot i) \quad \text{for } i \in \{3, 4\}$$

$$\text{AT3} = \neg(\square_{[10,30]} v \in [50, 60])$$

$$\text{AT4}(\bar{v}, \underline{\omega}) = (\square_{[0,10]} v < \bar{v}) \vee (\diamond_{[0,30]} \underline{\omega} < \omega)$$

There are a few noteworthy details. The robustness score of AT2 remains ∞ until $g = i$ is found by unguided search, and will jump abruptly to a finite value at that time. To falsify AT3, a fairly precise speed between 50 and 60 has to be reached and maintained, which requires a fine grained control of the input, which could be a challenge for aLVTS. To falsify AT4 one must balance keeping the engine RPM $\omega \leq \underline{\omega}$ while accelerating sufficiently to reach speed $v > \bar{v}$ before the deadline at time 10.

The input signal for the benchmarks is piecewise constant with 4 control points for random sampling/Breach, which is sufficient to falsify all requirements. We choose 6 levels, with 2,2,3,3,3,4 control points, respectively, corresponding to a time granularity of input segment durations between 15 ($= \frac{30}{2}$, coarsest) to 7.5 ($= \frac{30}{4}$, finest).

Powertrain Control. The benchmark was proposed for hybrid systems verification and falsification in [22]. It comprises a detailed model of a fuel control system that maintains the air-to-fuel AF ratio close to a reference value AF_{ref} . The model has two control algorithms: During startup and when high power is requested, it operates in feed-forward mode, whereas otherwise measurements of the actual air-to-fuel are fed back into the controller. Falsification tries to detect amplitude and duration of spikes in the air-to-fuel ratio that occur as a response either to mode switches or to changes in the throttle θ at a given engine speed ω . The relation between

the three quantities is nontrivial and in part determined by time delays. The input $\theta \in [0, 62.1)$ (normal mode) and $\theta \in [61.2, 81.2)$ (power mode) varies throughout the trace, whereas $\omega \in [900, 1100]$ is initially chosen and kept constant. Requirements are expressed with respect to the normalized error $\mu = |AF - AF_{\text{ref}}|/AF_{\text{ref}}$, where AF_{ref} depends dynamically on the current mode.

We consider requirement 27 from [22], denoted AFC27 in the following, which states that after a falling or rising edge of the throttle, μ should return to the reference value within 1 time unit and stay close to it for some time. We represent this as

$$\text{AFC27} = \square_{[11,50]} (\text{rise} \vee \text{fall}) \implies \square_{[1,5]} |\mu| < 0.008,$$

where rising and falling edges of θ are detected by $\text{rise} = \theta < 8.8 \wedge \diamond_{[0,\epsilon]} 40.0 < \theta$ and $\text{fall} = 40.0 < \theta \wedge \diamond_{[0,\epsilon]} \theta < 8.8$ for $\epsilon = 0.1$. The concrete bound of 0.008 in AFC27 and the edge detection parameters are taken from the report of the ARCH friendly competition on falsification (which featured this problem in 2017 [9] and 2018 [10]) as a balance between difficulty of the problem and ability to find falsifying traces. The interval [11, 50] over which the requirement is checked begins 1 time unit after the transition from startup to normal mode at time 10. Without this short delay, falsification would trivially discover the large spike following the mode switch.

The input signal is piecewise constant with 10 control points, again following [9, 10], and 5 levels with 10 control points each for aLVTS.

4.3 Experimental Results

We compare the performance of aLVTS with uniform random sampling (both implemented in FALSTAR²) and with the state-of-the-art stochastic global optimization algorithm CMA-ES [19] implemented in the falsification tool Breach³. We do not make a comparison with Breach's Nelder-Mead algorithm, since it has significantly poorer performance than CMA-ES on almost all of the benchmarks. The machine and software configuration was: CPU Intel i7-3770, 3.40GHz, 8 cores, 8Gb RAM, 64-bit Ubuntu 16.04 kernel 4.4.0, MATLAB R2018a, Scala 2.12.6, Java 1.8.

We compare two performance metrics: *success rate* (how many falsification trials were successful in finding a falsifying input) and the *number of iterations* made, which corresponds to the number of simulations required and thus indicates time. To account for the stochastic nature of the algorithms, the experiments were repeated for 50 trials. For a meaningful comparison of the number of iterations until falsification, we tried to maximize the falsification rate for a limit of 300 iterations per trial.

The number of iterations of the top-level loop in algorithm 1 in our implementation corresponds exactly to one complete Simulink simulation up to the time horizon (cf. section 4.1). For random sampling and CMA-ES, the number of iterations likewise corresponds to samples taken by running exactly one simulations each. Hence the comparison is fair and, as the overhead is dominated by simulation time, the numbers are roughly proportional to wall-clock times.

Table 2 summarizes our results in terms of success rate, and average number of iterations (M) and standard deviation (SD) of

²<https://github.com/ERATOMMSD/falstar> commit 43f5ca7

³<https://github.com/decyphir/breach> release version 1.2.9

Table 2: Experimental Results comparing the success rate out of 50 trials (“succ.”, higher is better) and number of iterations averaged over successful trials (“iter.”, lower is better, where M=arithmetic mean, SD=standard deviation) of uniform random sampling, Breach/CMA-ES, and FALSTAR/algorithm 1 for a maximum of 300 iterations per trial. The best results for each requirement are highlighted.

Formula	Random			Breach: CMA-ES			FALSTAR: aLVTS		
	succ. /50	iter. M	SD	succ. /50	iter. M	SD	succ. /50	iter. M	SD
AT1	43	106.6	83.9	50	39.7	23.6	50	8.5	6.7
AT2 ($i = 3$)	50	41.0	36.7	50	13.2	9.1	50	33.4	27.5
AT2 ($i = 4$)	49	67.0	60.8	6	17.8	15.9	50	23.4	22.5
AT3	19	151.1	98.1	50	145.2	63.0	50	86.3	52.1
AT4 (a)	36	117.3	71.8	50	97.0	47.7	50	22.8	10.6
AT4 (b)	2	117.7	9.2	49	46.7	58.0	50	47.6	23.5
Summary AT	199	95.3	47.9	255	42.8	29.0	300	29.2	19.4
AFC27	15	129.1	90.8	41	121.0	49.3	50	3.9	4.3

Parameters AT4 (a): $\bar{v} = 80$, $\underline{\omega} = 4500$, AT4 (b): $\bar{v} = 50$, $\underline{\omega} = 2700$.
Summary AT for succ.: sum, for iter.: geometric mean (see [17]).

successful trials. The unambiguously (possibly equal) best results are highlighted in blue. Where the lowest average number of iterations was achieved without finding a falsifying input for every trial, we highlight in grey the lowest average number of iterations for 100% success. We thus observe that aLVTS achieves the best performance in all but one case, AT2 ($i = 3$). Importantly, within the budget of 300 iterations per trial, aLVTS achieves a perfect success rate. CMA-ES is successful in 296 trials out of the total 350, with sub maximal success for AT2 ($i = 4$) and AFC27. In comparison, random sampling succeeds in only 214 trials, with sub maximal success in all but AT2.

The number of iterations required for falsification varies significantly between the algorithms and between the benchmarks. For the automatic transmission benchmarks, as an approximate indication of relative performance, CMA-ES requires about 50% more iterations as aLVTS (geometric mean: 42.8 vs. 29.2), and random sampling requires again twice as many as CMA-ES. Except for AT2, variance of aLVTS is lower than that of CMA-ES.

Looking at the individual benchmarks, CMA-ES is only consistently faster than aLVTS on AT2 ($i = 3$). Despite having fewer average iterations on AT2 ($i = 4$), CMA-ES cannot be considered faster than aLVTS because most of its trials fail to find a falsifying input—the numbers of iterations are not comparable. For the powertrain model (AFC27), the performance of aLVTS is more than an order of magnitude better: 3.9 iterations on average, compared to 121 for CMA-ES.

Figure 2 compares all trial runs for AFC27, ordered by the number of iterations required for falsification. Similar plots for the automatic transmission benchmarks are shown in fig. 3. The shape of each curve gives an intuition of the performance and consistency of its corresponding algorithm. In general, fewer iterations and more

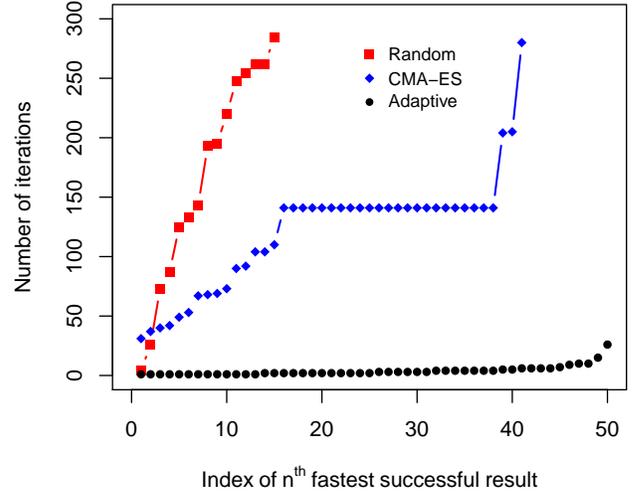


Figure 2: Relative performance of falsification algorithms with AFC27 powertrain requirement: fewer iterations and more successful results are better.

successful results are better, so it is immediately clear that our aLVTS performs significantly better than both random sampling and CMA-ES.

4.4 Discussion

For AT1, aLVTS quickly finds the falsifying input signal, as the required throttle of 100 and brake of 0 are contained in level 0 and are very likely to be tried early on. In contrast, even though this is a problem that is well-suited to hill-climbing, CMA-ES has some overhead to sample its initial population, as clearly visible in fig. 4a.

While CMA-ES deals very well with AT2 for $i = 3$, it struggles to find falsifying inputs for $i = 4$ (cf. figs. 4b and 4c). We attribute this to the fact that reaching gear 4 by chance occurs rarely in the exploration of CMA-ES when the robustness score is uninformative. aLVTS not only explores the spatial dimensions, but takes opportunistic jumps to later time points, which increases the probability of discovering a trace (prefix) where the gear is reached.

A priori, one would expect CMA-ES to perform well with AT3 and AT4, exploiting its continuous optimization to fine tune inputs between conflicting requirements. E.g., AT3 requires that v is both above 50 and below 60; AT4 requires that v is high while maintaining low ω , which is proportional to v . One would similarly expect the limited discrete choices made by aLVTS to hinder its ability to find falsifying inputs. Despite these expectations, our results demonstrate that in most situations aLVTS converges to a falsifying input more consistently and with fewer iterations than CMA-ES. We speculate that this is because CMA-ES is too slow to reach the “sweet spots” in the input space, where its optimization is efficient.

For AT3, there are a few instances where where aLVTS does not quickly find a good prefix towards the corridor $v \in [50, 60]$ at time 10 (the rightmost points in fig. 4d), which can be explained by the probabilistic nature of the search.

Regarding two instances of AT4 in figs. 4e and 4f, the graph for aLVTS is generally smoother and shallower, whereas CMA-ES shows consistent performance for only some of the trials but takes significantly more time on the worst 10 trials. We remark that the two parameter settings seem to pose opposite difficulty for the two algorithms, as CMA-ES is significantly quicker for two thirds of the trials for the second instance. It is unclear what causes this variation in performance of CMA-ES.

The plateaux apparent in some of the results using CMA-ES are difficult to explain, but suggest some kind of procedural logic or counter to decide termination. In contrast, the curves for random sampling and aLVTS are relatively smooth, reflecting their purely probabilistic natures.

5 RELATED WORK

The idea to find falsifying inputs using robustness as an optimization function for global optimization originates from [16] and has since been extended to the parameter synthesis problem (e.g., [23]). Two mature implementations in MATLAB are S-Taliro [4] and Breach [11], which have come to define the benchmark in this field. We describe their optimization algorithms below. Approaches to make the robustness semantics more informative are [3, 15]. These use integrals instead of min/max in the semantics of temporal operators.

Underminer [5] is a recent falsification tool that learns the (non-) convergence of a system to direct falsification and parameter mining. It supports STL formulas, SVMs, neural nets, and Lyapunov-like functions as classifiers. Other global approaches include [1], which partitions the input space into sub-regions from which falsification trials are run selectively. This method uses coverage metrics to balance exploration and exploitation. Comprehensive surveys of simulation based methods for the analysis of hybrid systems are given in [7, 24].

Users of S-Taliro and Breach can select from a range of optimization algorithms (Uniform Random, Nelder-Mead, Simulated Annealing, Cross-Entropy, CMA-ES). These cover a variety of trade-offs between exploration of the search space and exploitation of known good intermediate results. As such, their performance varies with the structure of the problem at hand.

In general, global optimization needs to find a solution in an *unstructured* combinatorial search space, where the different parameters are assumed, a priori, to be independent. Good combinations of parameter choices have to be discovered first, before the optimization can work effectively (cf. AT2 and AT3). This incurs an exponential cost of exhausting the different combinations, which is mitigated in different ways by the various algorithms. Nelder-Mead is seeded by a constant number of samples and tries to obtain better points in the search space by making linear combinations of the previous ones, according to a fixed scheme. The price paid is not just lack of exploration, but also the danger of selecting mid points between good and *bad* values in some dimensions, as a consequence of the weighted average of linear combinations. To this end, Cross-Entropy methods and CMA-ES discover relations between parameters and exploit these by decorrelating the choices across independent dimensions (i.e., not trying all possible combinations), as explained nicely in the context of falsification in [28]. Simulated

Annealing, on the other hand, relies on fairly unguided random exploration, and thus requires many iterations to converge, which can be prohibitively expensive for complex models.

In comparison, our approach explicitly takes into account the temporal structure of the falsification problem. The exponential search space is pruned by focusing on the early part of the input signal, by investigating a small number of coarse choices first, and by ignoring those prefixes which are unpromising according to the robustness heuristic. As demonstrated by our experiments, the strength of optimization algorithms to discover precise solutions in the continuous input space does not seem efficient when the mass of falsifying traces is large and exploration is the key factor.

The characteristic of our approach to explore the search space incrementally is shared with *rapidly-exploring random trees* (RRTs). The so-called star discrepancy metric guides the search towards unexplored regions and a local planner extends the tree at an existing node with a trajectory segment that closely reaches the target point. RRTs have been used successfully in robotics [26] and also in falsification [14]. On the other hand, the characteristic of taking opportunistic coarse jumps in time is reminiscent of stochastic local search [8] and multiple-shooting [33].

Monte Carlo tree search (MCTS) has been applied to a model of aircraft collisions in [27]; and in a falsification context more recently [32] to guide global optimization, building on the previous idea of time-staging [31]. That work noted the strong similarities between falsification using MCTS and statistical model checking (SMC) using *importance splitting* [20]. The robustness semantics of STL, used in [31, 32] and the present approach to guide exploration, can be seen as a “heuristic score function” [21] in the context of importance splitting. All these approaches construct trees from traces that share common prefixes deemed good according to some heuristic. The principal difference is that importance splitting aims to construct a diverse set of randomly-generated traces that all satisfy a property (equivalently, falsify a negated property), while falsification seeks a single falsifying input. The current work can be distinguished from standard MCTS and reinforcement learning [29] for similar reasons. These techniques tend to seek optimal policies that make good decisions in *all* situations, unnecessarily covering the entire search space.

6 CONCLUSION

We have presented a probabilistic algorithm that finds inputs to a hybrid system that falsify a given temporal logic requirement. On standard benchmarks our algorithm significantly and consistently outperforms existing methods. While the falsification problem is inherently hard (no theoretically best solution can exist), we believe that we have found an approach that gives good results in practice, with probabilistic certainty. We believe the reason it works well stems from a property shared by many falsification problems, namely, that there tends to be a significant mass of falsifying inputs. By design it scales with the difficulty of the problem, finding trivial solutions (extreme inputs) immediately and has a small constant overhead for “almost easy” problems. This pays off, in particular, for system models that are expensive to simulate, such as the powertrain benchmark. The approach admits incremental computation of simulations and can easily be parallelized.

As future work, we wish to extend the method with some form of optimization, e.g., by interpolation of previously seen traces, reminiscent of the linear combinations computed by Nelder-Mead. We expect that this would help to “fine-tune” results with small positive robustness into full solutions, mitigating the limitation inherent in the discretization of the input space. Likewise, extrapolation of results could be used to propagate the robustness from one level to the next, reducing the need to look at different branches with similar inputs. Finally, we plan to apply the approach to more benchmarks, specifically those that combine discrete and continuous input domains.

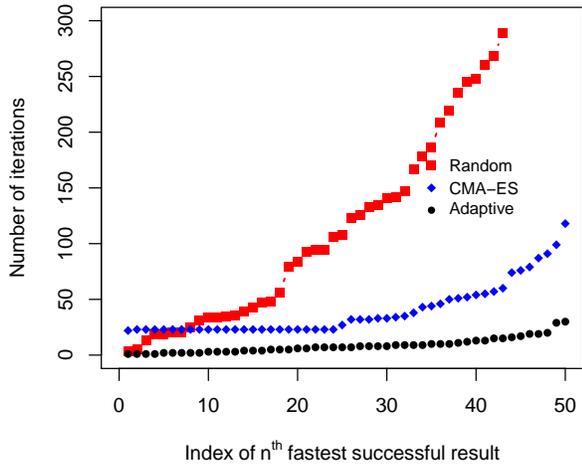
ACKNOWLEDGMENTS

The authors are supported by ERATO HASUO Metamathematics for Systems Design Project (No. JPMJER1603), Japan Science and Technology Agency; and Grants-in-Aid No. 15KT0012, Japan Society for the Promotion of Science.

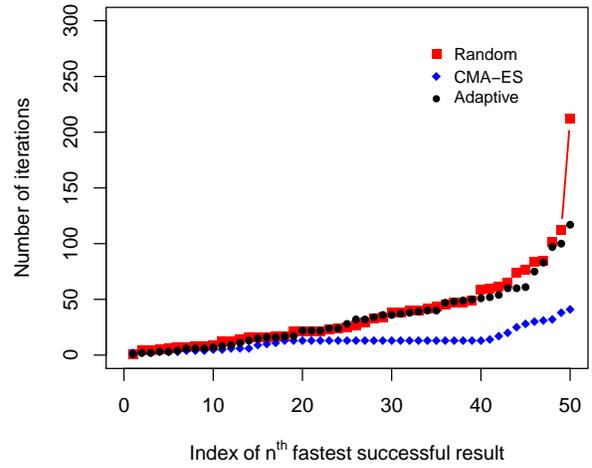
REFERENCES

- [1] Arvind Adimoolam, Thao Dang, Alexandre Donzé, James Kapinski, and Xiaoqing Jin. 2017. Classification and Coverage-Based Falsification for Embedded Control Systems. In *Computer Aided Verification*, Rupak Majumdar and Viktor Kunčák (Eds.). Springer, Cham, 483–503.
- [2] Takumi Akazaki. 2016. Falsification of Conditional Safety Properties for Cyber-Physical Systems with Gaussian Process Regression. In *Runtime Verification (RV)*, Yliès Falcone and César Sánchez (Eds.). LNCS, Vol. 10012. Springer, 439–446.
- [3] Takumi Akazaki and Ichiro Hasuo. 2015. Time robustness in MTL and expressivity in hybrid system falsification. In *Computer Aided Verification (CAV) (LNCS)*, Vol. 9207. Springer, 356–374.
- [4] Yashwanth Annpureddy, Che Liu, Georgios Fainekos, and Sriram Sankaranarayanan. 2011. S-TaLiRo: A Tool for Temporal Logic Falsification for Hybrid Systems. In *Tools and Algorithms for the Construction and Analysis of Systems (TACAS) (LNCS)*, Parosh Aziz Abdulla and K. Rustan M. Leino (Eds.). Springer, 254–257.
- [5] Ayca Balkan, Paulo Tabuada, Jyotirmoy V. Deshmukh, Xiaoqing Jin, and James Kapinski. 2017. Underminer: A Framework for Automatically Identifying Non-converging Behaviors in Black-Box System Models. *ACM Trans. Embed. Comput. Syst.* 17, 1, Article 20 (2017), 28 pages.
- [6] Clark Barrett, Pascal Fontaine, and Cesare Tinelli. 2017. *The SMT-LIB Standard: Version 2.6*. Technical Report. Department of Computer Science, The University of Iowa. Available at www.SMT-LIB.org.
- [7] Ezio Bartocci, Jyotirmoy Deshmukh, Alexandre Donzé, Georgios Fainekos, Oded Maler, Dejan Ničković, and Sriram Sankaranarayanan. 2018. Specification-based monitoring of cyber-physical systems: a survey on theory, tools and applications. In *Lectures on Runtime Verification*. Springer, 135–175.
- [8] Jyotirmoy Deshmukhand, Xiaoqing Jin, James Kapinski, and Oded Maler. [n. d.]. Stochastic Local Search for Falsification of Hybrid Systems. *Automated Technology for Verification and Analysis (ATVA) 9364* ([n. d.]), 500–517.
- [9] Adel Dokhanchi, Shakiba Yaghoubi, Bardh Hoxha, and Georgios E. Fainekos. 2017. ARCH-COMP17 Category Report: Preliminary Results on the Falsification Benchmarks. In *Applied Verification of Continuous and Hybrid Systems (ARCH) (EPIc Series in Computing)*, Goran Frehse and Matthias Althoff (Eds.), Vol. 48. EasyChair, 170–174.
- [10] Adel Dokhanchi, Shakiba Yaghoubi, Bardh Hoxha, Georgios E. Fainekos, Gidon Ernst, Zhenya Zhang, Paolo Arcaini, Ichiro Hasuo, and Sean Sedwards. 2019. ARCH-COMP18 Category Report: Results on the Falsification Benchmarks. In *Applied Verification of Continuous and Hybrid Systems (ARCH) (EPIc Series in Computing)*, Goran Frehse (Ed.), Vol. 54. EasyChair, 104–109.
- [11] Alexandre Donzé. 2010. Breach, A toolbox for verification and parameter synthesis of hybrid systems. In *International Conference on Computer Aided Verification (CAV 2010) (LNCS)*. Springer, 167–170.
- [12] Alexandre Donzé, Thomas Ferrère, and Oded Maler. 2013. Efficient Robust Monitoring for STL. In *Computer Aided Verification (CAV) (LNCS)*, Natasha Sharygina and Helmut Veith (Eds.), Vol. 8044. Springer, 264–279.
- [13] Alexandre Donzé and Oded Maler. 2010. Robust Satisfaction of Temporal Logic over Real-Valued Signals. In *Formal Modeling and Analysis of Timed Systems (FORMATS 2010) (LNCS)*, Krishnendu Chatterjee and Thomas A. Henzinger (Eds.). Springer, 92–106.
- [14] Tommaso Dreossi, Thao Dang, Alexandre Donzé, James Kapinski, Xiaoqing Jin, and Jyotirmoy V. Deshmukh. 2015. Efficient Guiding Strategies for Testing of Temporal Properties of Hybrid Systems. In *NASA Formal Methods (NFM) (LNCS)*, Klaus Havelund, Gerard Holzmann, and Rajeev Joshi (Eds.), Vol. 9058. Springer, 127–142.
- [15] Johan Eddeland, Sajed Miremadi, Martin Fabian, and Knut Åkesson. 2017. Objective functions for falsification of signal temporal logic properties in cyber-physical systems. In *Conference on Automation Science and Engineering (CASE)*. IEEE, 1326–1331.
- [16] Georgios E. Fainekos and George J. Pappas. 2009. Robustness of temporal logic specifications for continuous-time signals. *Theor. Comp. Sci.* 410, 42 (2009), 4262–4291.
- [17] Philip J Fleming and John J Wallace. 1986. How not to lie with statistics: the correct way to summarize benchmark results. *Commun. ACM* 29, 3 (1986), 218–221.
- [18] Bardh Hoxha, Houssam Abbas, and Georgios E. Fainekos. 2014. Benchmarks for Temporal Logic Requirements for Automotive Systems. In *Applied Verification for Continuous and Hybrid Systems (ARCH) (EPIc Series in Computing)*, Goran Frehse and Matthias Althoff (Eds.), Vol. 34. EasyChair, 25–30.
- [19] Christian Igel, Nikolaus Hansen, and Stefan Roth. 2007. Covariance matrix adaptation for multi-objective optimization. *Evolutionary computation* 15, 1 (2007), 1–28.
- [20] Cyrille Jegourel, Axel Legay, and Sean Sedwards. 2013. Importance Splitting for Statistical Model Checking Rare Properties. In *Computer Aided Verification (CAV)*. LNCS, Vol. 8044. Springer, 576–591.
- [21] Cyrille Jegourel, Axel Legay, and Sean Sedwards. 2014. An Effective Heuristic for Adaptive Importance Splitting in Statistical Model Checking. In *Leveraging Applications of Formal Methods, Verification and Validation. Specialized Techniques and Applications (SoLA)*, Tiziana Margaria and Bernhard Steffen (Eds.). LNCS, Vol. 8803. Springer, 143–159.
- [22] Xiaoqing Jin, Jyotirmoy V. Deshmukh, James Kapinski, Koichi Ueda, and Kenneth R. Butts. 2014. Powertrain control verification benchmark. In *Hybrid Systems: Computation and Control (HSCC)*, Martin Fränzle and John Lygeros (Eds.). ACM, 253–262.
- [23] Xiaoqing Jin, Alexandre Donzé, Jyotirmoy V. Deshmukh, and Sanjit A. Seshia. 2015. Mining requirements from closed-loop control models. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 34, 11 (2015), 1704–1717.
- [24] James Kapinski, Jyotirmoy V. Deshmukh, Xiaoqing Jin, Hisahiro Ito, and Ken Butts. 2016. Simulation-Based Approaches for Verification of Embedded Control Systems: An Overview of Traditional and Advanced Modeling, Testing, and Verification Techniques. *IEEE Control Systems Magazine* 36, 6 (Dec 2016), 45–64.
- [25] Levente Kocsis and Csaba Szepesvári. 2006. Bandit Based Monte-Carlo Planning. In *Machine Learning (ECML 2006) (LNCS)*, Johannes Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou (Eds.), Vol. 4212. Springer, 282–293.
- [26] Steven M LaValle and James J Kuffner Jr. 2001. Randomized kinodynamic planning. *The International Journal of Robotics Research (IJRR)* 20, 5 (2001), 378–400.
- [27] Ritchie Lee, Mykel J. Kochenderfer, Ole J. Mengshoel, Guillaume P. Brat, and Michael P. Owen. 2015. Adaptive stress testing of airborne collision avoidance systems. In *IEEE/AIAA 34th Digital Avionics Systems Conference (DASC 2015)*. 6C2:1–6C2:13.
- [28] Sriram Sankaranarayanan and Georgios Fainekos. 2012. Falsification of temporal properties of hybrid systems using the cross-entropy method. In *Hybrid Systems: Computation and Control (HSCC)*. ACM, 125–134.
- [29] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction* (second ed.). MIT press.
- [30] David Wolpert and William G. Macready. 1997. No free lunch theorems for optimization. *IEEE Trans. Evolutionary Computation* 1, 1 (1997), 67–82.
- [31] Zhenya Zhang, Gidon Ernst, Ichiro Hasuo, and Sean Sedwards. 2018. Time-Staging Enhancement of Hybrid System Falsification. In *2018 IEEE Workshop on Monitoring and Testing of Cyber-Physical Systems (MT-CPS 2018)*. IEEE, 3–4.
- [32] Zhenya Zhang, Gidon Ernst, Sean Sedwards, Paolo Arcaini, and Ichiro Hasuo. 2018. Two-Layered Falsification of Hybrid Systems Guided by Monte Carlo Tree Search. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD 2018)* (2018).
- [33] Aditya Zutshi, Jyotirmoy V. Deshmukh, Sriram Sankaranarayanan, and James Kapinski. 2014. Multiple shooting, CEGAR-based falsification for hybrid systems. In *Embedded Software (EMSOFT)*. 5:1–5:10.

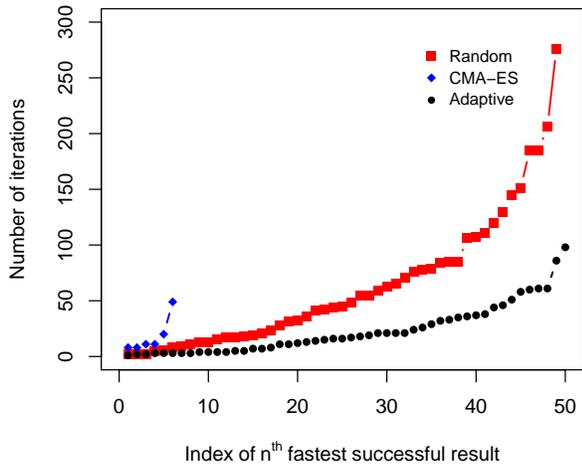
Figure 3: Performance comparison for the automatic transmission benchmark



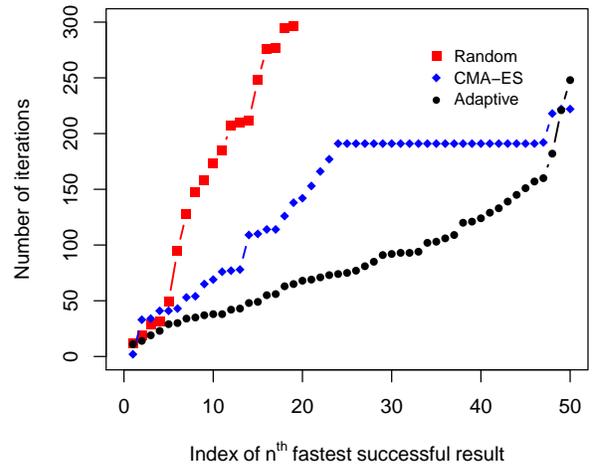
(a) Performance plot for AT1



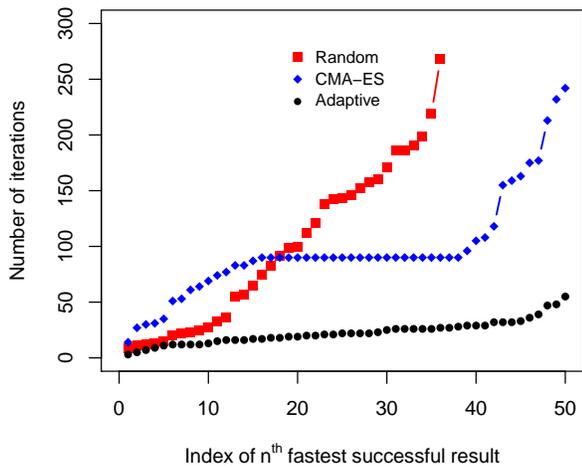
(b) Performance plot for AT2 ($i = 3$)



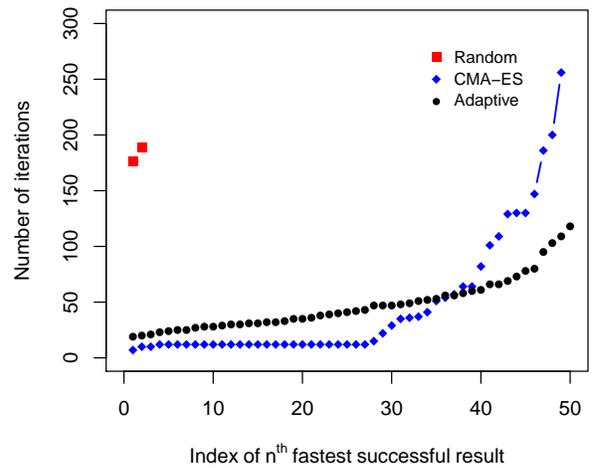
(c) Performance plot for AT2 ($i = 4$)



(d) Performance plot for AT3



(e) Performance plot for AT4 ($\bar{v} = 80, \underline{\omega} = 4500$)



(f) Performance plot for AT4 ($\bar{v} = 50, \underline{\omega} = 2700$)