

# Classification and coverage-based falsification for embedded control systems

Arvind Adimoolam<sup>1</sup>, Thao Dang<sup>1</sup>, Alexandre Donzé<sup>2</sup>,  
James Kapinski<sup>3</sup>, Xiaoqing Jin<sup>3</sup>

<sup>1</sup> CNRS/Verimag, Grenoble, France `santosh.adimoolam@univ-grenoble-alpes.fr`, `thao.dang@univ-grenoble-alpes.fr`

<sup>2</sup> Decyphir, Inc, San Francisco, CA `alex@decyphir.com`

<sup>3</sup> Toyota Motors North America R&D, Gardena, CA  
`{firstname.lastname}@toyota.com`

**Abstract.** Many industrial cyber-physical system (CPS) designs are too complex to formally verify system-level properties. A practical approach for testing and debugging these system designs is falsification, wherein the user provides a temporal logic specification of correct system behaviors, and some technique for selecting test cases is used to identify behaviors that demonstrate that the specification does not hold for the system. While coverage metrics are often used to measure the exhaustiveness of this kind of testing approach for software systems, existing falsification approaches for CPS designs do not consider coverage for the signal variables. We present a new coverage measure for continuous signals and a new falsification technique that leverages the measure to efficiently identify falsifying traces. This falsification algorithm combines global and local search methods and uses a classification technique based on support vector machines to identify regions of the search space to focus exploration effort. We use an industrial example from an automotive fuel cell application and other benchmark models to compare the new approach against some existing falsification tools.

## 1 Introduction

Cyber-physical systems integrate heterogeneous components whose descriptions in high level modeling languages involve a wide array of specification paradigms, such as differential equations, difference equations, automata, and data flow graphs. Although the behavior of individual cyber-physical components may be amenable to rigorous mathematical reasoning and analysis, the complex interactions between the components are still not well-understood and pose major theoretical hurdles in formal reasoning. Also, the scalability of the existing formal verification methods and tools (see [16, 3, 18, 4, 26, 1, 12, 15] and references therein) is still limited and therefore not suited for verifying industrial scale cyber-physical systems. Testing is an alternative approach for detecting errors, whose advantage over formal verification methods is that it can treat a system as a black box, meaning that no internal description of the system is required. In black box testing, only an interface of the system with the external environment is described. Although testing can be applied to large scale cyber-physical systems, as attested by its use in industry, it does not provide proofs of correctness. In other words, black box

testing can only detect bugs, and when it does so successfully, the system design has to be corrected. Nevertheless, when the testing process does not find any bugs, we cannot draw any conclusion about its correctness. If the falsification is unsuccessful, then information about the potential validity of the correct behavior of the system would be of great interest to the designer. This information can be provided in terms of a testing coverage measure.

In the existing research on cyber-physical systems testing, the focus was generally on *state-coverage measures*, that is measures to characterize the portion of the state space covered by a test suite. An example is star discrepancy [6, 11], a notion borrowed from statistics that indicates how equi-distributed are a set of tested points in the state space. Some other measures are dispersion [13], which indicates the size of the largest unexplored areas, and grid-cell count [25]. Although these state-coverage measures can serve as a possible means to compare coverage of testing data generated by different algorithms, these measures exhibit the following drawbacks. Typically, a test generation algorithm guided by a state coverage measure tries to sample test cases in the areas that are not well explored; however, in industrial models describing interactions among a large number of heterogeneous components, information about the state can be hard to obtain. Additionally, such systems can have low controllability, meaning that it is difficult or impossible to reach some regions of the state space. In such a case, the algorithms can expend a large amount of time attempting to explore unreachable regions. So, state-coverage measures are not appropriate for analysis or guidance of testing effort on many cyber-physical systems.

The present work addresses the shortcomings of the state-coverage based techniques by instead focusing on *coverage of input signal spaces*. We develop a new test generation technique that is based on covering the input signal space rather than the state space. Note that previous test generation methods have considered coverage of a parameter space (such as in [9]); the way that we handle input signals is directly related to them, as we consider the class of finitely parameterized input signals.

While coverage is important for providing confidence in correctness of the system behavior, bug detection is still an important goal of testing. Usually, there is a mutual tradeoff between satisfying the two criteria. Achieving good coverage entails exploring a large portion of the search space, most of which would correspond to correct behaviors. Whereas, the objective of a falsification procedure is to find incorrect behaviors, which would require focusing on behaviors close to incorrectness. Most falsification methods are based on minimizing the behavioral robustness with respect to a property under test; the robustness measure here indicates how far the behavior is from violating the property. A common drawback of such falsification methods is that the optimization procedures can spend a significant amount of computation time near local optima that may not correspond to an incorrect behavior. Therefore, a criterion like coverage can help overcome this drawback, since seeking to improve the global coverage would drive the search process out of the areas of local optima. One way to achieve a good compromise between coverage-driven and local search-driven testing is to initiate the search procedures from points that are separated by some threshold distance. This insight was used previously in the tabu search method, which ensures that all the starting points are well separated [7]; however, apart from ensuring that the starting points are well separated,

it is also desirable that they are chosen in regions in which one can reasonably expect to find an incorrect behavior. That is, heuristically speaking, a starting point should have a low robustness value.

Based on the above observations, in this work we present a falsification algorithm that combines the following three essential ideas:

- Defining a coverage measure for quantifying the exploration of input signal space during testing.
- Guiding a randomized global search procedure by performing robustness classification: the classification divides the search space into regions with different potentials of falsification characterized by the robustness of evaluated test cases. Our classification is inspired from linear *support vector machines* [8, 17].
- Using local search in regions classified as less robust. The above-mentioned global search together with an iterative classification procedure does converge towards an incorrect behavior, if it exists. However, to speed up the convergence, instead of continually classifying, we can use the information obtained from classification to efficiently initialize a local search within each classified region. Note that in general, local search with arbitrary initialization can perform poorly. Therefore, by alternating classification and local search we can achieve a better convergence while assuring a good coverage of the input signal space (because in general local search does not take into account this coverage criterion).

For implementation and evaluation purposes, we use a local search method, called the CMA-ES (Covariance Matrix Adaptation Evolution Strategy) [21], also used by the tool Breach [9]. The CMA-ES algorithm is considered as the state-of-the-art in evolutionary computation and has been used for industrial optimization applications. It is nevertheless important to note that any existing local search method can be used in our approach. The experimental results obtained using a MATLAB implementation of our falsification algorithm on some benchmark systems demonstrate its good performance and, in addition, its efficiency improvements over search algorithms like the CMA-ES. Indeed, our algorithm was tested on a difficult property of the PTC benchmark [11] and could falsify in all the tested random seeds while the methods based on pseudo-random sampling or only on the CMA-ES could not. Also, we demonstrate that the technique can be successfully applied to industrial problems by presenting results for a prototype automotive hydrogen fuel cell application.

Our approach draws inspiration from the approaches implemented in the tools S-Taliro [2] and Breach [9]. These approaches seek the worst case behaviors using the notion of robustness metrics, which are defined with respect to properties specified using the languages MTL (Metric Temporal Logic) [14] and STL (Signal Temporal Logic) [10]. The tools identify property violations by employing global optimization methods to search for behaviors that minimize robustness, where negative robustness values correspond to property violations. Robustness-based approaches can be seen as complementary to coverage-based approaches, since the former try to find a worst-case behavior while the latter try to cover a large number of possible behaviors. When a robustness-based approach cannot find an erroneous behavior due to the limitations of global optimization algorithms, the observed error absence cannot be used as a formal correctness proof; in this case good coverage would be desirable to enhance the

confidence that the system is error-free. By combining robustness-based and coverage-guided explorations, our approach enhances the overall testing effectiveness by providing confidence that important or representative behaviors are tested. Finally, we remark that the idea of combining global and local search was investigated in work [22], however this work requires computing sensitivity of state trajectories with respect to parameters and initial conditions, which is often impossible for black-box systems.

## 2 Preliminaries

We consider system models defined by a mapping from parameters and input signals to output signals,

$$y = \Phi(v, u), \quad (1)$$

where  $v \in \mathcal{V}$  is a valuation of a finite collection of parameters, and  $u \in \mathcal{U}$  is an input signal used to simulate the system. In this setting,  $\mathcal{V}$  could contain a set of system initial conditions as well as some finite set of system parameters. Each input signal  $u \in \mathcal{U}$  is a function  $\mathcal{I}_u \mapsto U$ , where  $\mathcal{I}_u$  is an interval (either discrete or continuous) from 0 to some finite value, and  $U$  is some metric space of finite dimension. Similarly, we assume that each output signal  $y \in \mathcal{Y}$  is a function  $\mathcal{I}_y \mapsto Y$ , where  $\mathcal{I}_y$  is an interval (either discrete or continuous) from 0 to some finite value, and  $Y$  is some metric space of finite dimension. We assume that  $\mathcal{V}, \mathcal{U}$ , and  $\mathcal{Y}$  are metric spaces. Note that the system defined by (1) does not explicitly model the behaviors of the internal system states. State behaviors could be modeled using this framework by ensuring that  $v$  also includes the system state and all of the states map to system outputs, but we do not require this.

We assume that signals are finitely parameterized, i.e., an input signal  $u$  can be uniquely determined by a finite set of  $m$  parameters, whose valuation  $\hat{u}$  is in a subset  $\hat{\mathcal{U}}$  of an  $m$ -dimensional metric space. For example, a right-continuous piecewise constant input signal  $u : \mathcal{I}_u \rightarrow \mathbb{R}$ , where  $\mathcal{I}_u = [0, T]$ , with discontinuities occurring at monotonically increasing instants  $\tau_1, \dots, \tau_m$ , where  $0 = \tau_1 < \tau_m < T$ , can be uniquely defined by  $m$  values  $u(\tau_i)$ . Subsequently, our system can be defined as a map from a finite set of parameters to the output signals, as follows:

$$y = \hat{\Phi}(v, \hat{u}), \quad v \in \mathcal{V} \text{ and } \hat{u} \in \hat{\mathcal{U}} \quad (2)$$

We call  $\mathcal{V}$  the space of *nominal parameters* and  $\hat{\mathcal{U}}$  the space of *input signal parameters*.

**Signal Temporal Logic.** To specify correct behavior of a system defined by (1), we use Signal Temporal Logic (STL) [23], which can capture behaviors of real-valued signals over discrete or dense time. We present here an informal description of STL (see [23] for more details). A formula in STL consists of atomic predicates, Boolean, and temporal operators. Atomic predicates over signal values are of the form  $\mu = f(y(t)) \sim 0$ , where  $f$  is a scalar-valued function over the signal  $y$  evaluated at time  $t$ , and  $\sim \in \{<, \leq, >, \geq, =, \neq\}$ . Temporal operators “always” ( $\Box$ ), “eventually” ( $\Diamond$ ), and “until” ( $\mathcal{U}$ ) have the usual meaning and are scoped using intervals of the form  $(a, b)$ ,  $(a, b]$ ,  $[a, b)$ ,  $[a, b]$ , or  $(a, \infty)$ , where  $a, b \in \mathbb{R}_{\geq 0}$  and  $a < b$ . If  $I$  is such an interval, then the language of STL is given by the following grammar:

$$\varphi := \top \mid f(y(t)) \sim 0 \mid \neg \varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \mathcal{U}_I \varphi_2 : \quad \sim \in \{<, \leq, >, \geq, =, \neq\} \quad (3)$$

The  $\Diamond$  and  $\Box$  operators are defined as follows:  $\Diamond_I \varphi \triangleq \top \mathcal{U}_I \varphi$ ,  $\Box_I \varphi \triangleq \neg(\Diamond_I \neg \varphi)$ . When omitted, the interval  $I$  is assumed to be  $[0, \infty)$ . The semantics are described informally as follows. The signal  $y$  satisfies  $f(y) > 0$  at time  $t$  if  $f(y(t)) > 0$ . It satisfies  $\varphi = \Box_{[0,1)}(f(y) = 0)$  if for all time  $0 \leq t < 1$ ,  $f(y(t)) = 0$ . The signal satisfies  $\varphi = \Diamond_{[1,2)}(f(y) < 0)$  iff there exists a time  $t$  such that  $1 \leq t < 2$  and  $f(y(t)) < 0$ . The two-dimensional signal  $y = (y_1, y_2)$  satisfies the formula  $\varphi = (y_1 > 0) \mathcal{U}_{[2.3, 4.5]}(y_2 < 0)$  iff there is some time  $t$  where  $2.3 \leq t \leq 4.5$ ,  $y_2(t) < 0$ , and  $\forall t'$  in  $[2.3, t)$ ,  $y_1(t') > 0$ .

**Quantitative Semantics for STL.** The quantitative semantics of STL tells how far a signal is from satisfying a formula. In this respect, we use the quantitative interpretation presented in [10], which we describe informally as follows. The semantics relies on a function  $\rho$  such that a positive sign of  $\rho(\varphi, y, t)$  indicates that  $(y, t)$  satisfies  $\varphi$ , and its absolute value estimates the *robustness* of this satisfaction. If  $\phi$  is a simple inequality of the form  $f(y) > b$ , then its robustness is  $\rho(\varphi, y, t) = f(y(t)) - b$ . For the conjunction of two formulas  $\varphi := \varphi_1 \wedge \varphi_2$ , we have  $\rho(\varphi, y, t) = \min(\rho(\varphi_1, y, t), \rho(\varphi_2, y, t))$ , while for the disjunction  $\varphi := \varphi_1 \vee \varphi_2$ , we have  $\rho(\varphi, y, t) = \max(\rho(\varphi_1, y, t), \rho(\varphi_2, y, t))$ . For a formula with until operator as  $\varphi := \varphi_1 \mathcal{U}_I \varphi_2$ , the robustness is computed as  $\rho(\varphi, y, t) = \max_{t' \in t+I} (\min(\rho(\varphi_2, y, t'), \min_{t'' \in [t, t']} (\rho(\varphi_1, y, t''))))$ .

Since the output signal is determined by the set of nominal parameters and input signal parameters according to the mapping  $\hat{\Phi}$ , we can define a robustness function over the space of parameters, called *parametric robustness*, as  $\hat{\rho}(\varphi, v, \hat{u}, t) = \rho(\varphi, \hat{\Phi}(v, \hat{u}), t)$ .

**Falsification.** Finding a counterexample of  $\varphi$  means finding a parameter value  $v \in \mathcal{V}$  and an input parameter value  $\hat{u} \in \hat{\mathcal{U}}$  such that  $y \not\models \varphi$ , where  $y = \hat{\Phi}(v, \hat{u})$ . Equivalently, the counterexample is identified when its parametric robustness is less than zero, i.e.,  $\hat{\rho}(\varphi, v, \hat{u}, t) < 0$  for some time point  $t$  in the time horizon of the signal. We call any  $v \in \mathcal{V}$  and  $\hat{u} \in \hat{\mathcal{U}}$  for which  $y \not\models \varphi$  a *counterexample* and we call this task of finding a counterexample as a *falsification problem*. We say that a counterexample  $y$  (that is  $y \not\models \varphi$ ) is *robust* if there exists a neighborhood around  $y$ ,  $\mathcal{N}_y$ , such that for all  $y' \in \mathcal{N}_y$ ,  $y' \not\models \varphi$ . We call a corresponding neighborhood  $\mathcal{N}_y$  a *robustness neighborhood* of counterexample  $y$ . If a counterexample has a robustness neighborhood that contains a closed ball of radius  $\epsilon$ , then we say that the counterexample is  $\epsilon$ -*robust*.

**Continuity of robustness.** Recall that our input signals are assumed to be finitely parametrized and correspondingly we defined the parametric robustness function. If we assume that the predicates of an STL formula are defined by functions  $f$  in (3) which are continuous w.r.t. the value of  $y$  at any time  $t$ , and the mapping  $\hat{\Phi}$  defining the system dynamics is continuous w.r.t. the parameter and input signal, then we can prove that the parametric robustness is continuous w.r.t.  $v$  and  $\hat{u}$ . Indeed, for any atomic predicate  $\varphi = f(y(t)) \sim 0$ , the parametric robustness  $\hat{\rho}(\varphi, v, \hat{u}, t)$  is continuous because  $f$  and  $\hat{\Phi}$  are continuous. Next, for any general formula as defined in (3), the robustness is computed by a composition of *min* and *max* operators of subformulas. By using induction we thus can deduce that the parametric robustness, given the aforementioned assumptions, is continuous in the input parameter  $\hat{u}$  and the nominal parameter  $v$ .

### 3 Input space coverage - Cell occupancy

This section presents a metric that we use to measure the coverage of signal spaces. The notion is intended to be used to define the coverage of input signals used to stimulate a dynamical system. We define a measure called cell occupancy, which has the following desirable properties:

- The measure is *monotonic*, in the sense that it is guaranteed not to decrease in value when new signals are added to an existing set;
- The measure permits computation with *efficient algorithms*;
- The measure provides numbers in *reasonable ranges*, in the sense that, for both low dimension and high dimension problems, the measure results in values that are neither too large nor too small so as to be accurately represented with standard floating point numbers.

Henceforth, we define a measure called *cell occupancy* as follows. Let  $M$  be a set of signals, which corresponds to a set of parameter vectors  $X_M$ . We call elements of  $X_M$  points.

Choose a partition of  $X$ ,  $\omega = \{\omega_i | i = 1, \dots, l\}$ . For now, we assume that each partition element, which we call a *cell*, is rectangular, with each side of equal length,  $\Delta$ , called *grid cell size*<sup>4</sup>. A vector that indicates how many points are in each cell is called a *distribution*,  $D = (n_1, \dots, n_l)$ , where each  $n_i$  indicates how many points are located in cell  $i$ . Cell occupancy is based on the relative number cells occupied by points, compared to the total number of cells. Consider the total number of occupied cells, that is, the number of cells that contain at least one point, i.e.,  $N_c = \sum_{i=1}^l g_i$  where  $g_i = 1$  if  $n_i \geq 1$ , and  $g_i = 0$  otherwise. Then, the proposed cell occupancy measure is given as

$$H_c(D) = \frac{\log N_c}{\log l}.$$

Logarithm functions are used due to the fact that the total number of cells could be very large as compared to the number of occupied cells. The logarithms provide two key features for the cell occupancy measure: (1) they maintain the monotonicity of the measure, and (2) they result in reasonable measure values even for cases where the dimension  $m$  is large.

*Guarantee for finding counterexample.* We consider here falsification algorithms based on an iterative search on the nominal parameter and input signal parameter spaces. We assume that the functions in the atomic predicates of STL formulas are continuous in the value of the output signal at any fixed time point. Also, the system mapping  $\hat{\phi}$  is assumed to be continuous w.r.t. the input signal parameters and nominal parameters.

<sup>4</sup> We note that in the setting in which we intend to apply the following coverage metrics, we will expect to select points in  $X$  that are no closer than some  $\epsilon$  distance from each other, based on some metric between signals, but this rectangularity will not be exploited in the following. Further, we assume that  $\epsilon \ll \Delta$ .

In this case, if a falsification algorithm is such that the cell occupancy is guaranteed to increase after a finite number of robustness evaluations for any partition, then because of the continuity of parametric robustness (explained in Section 2), there exists a sufficiently small upper bound on the grid cell size below which, the algorithm is guaranteed to find a counterexample. This is summarized by the following lemma. However, note that in general such falsification algorithms may be used for non-continuous systems as well. The following lemma gives a theoretical insight about why coverage may be taken into account for designing efficient falsification algorithms.

**Lemma 1.** *Given a falsification algorithm and a partition  $\omega$  with  $l$  grid cells of size  $\Delta > 0$ , let  $D(\kappa, \Delta)$  denote the cell distribution after  $\kappa$  robustness evaluations by the algorithm. Let us consider that there exists  $\alpha \in \mathbb{Z}_{>0}$  for which the algorithm guarantees that  $\forall \kappa \in \mathbb{Z}_{\geq 0} H_c(D(\kappa + \alpha, \Delta)) > H_c(D(\kappa, \Delta))$ . Let us also consider that the system mapping  $\hat{\Phi}$  is continuous and an STL formula  $\varphi$  is formed by continuous predicates with respect to the signal value at a fixed time. In this case, if an  $\epsilon$ -robust counterexample of  $\varphi$  exists, then there exists an upper bound  $\bar{\Delta} > 0$  on the grid cell size of  $\omega$  such that  $\forall \Delta < \bar{\Delta}$ , the algorithm finds a counterexample after a finite number of robustness evaluations.*

## 4 Falsification techniques

We use the term *sampling a point* to mean selecting a parameter vector  $x$  in the parameter set  $X_M$ , to uniquely define an input signal in  $M$ . Such signals are then used as stimuli to simulate the system and determine the robustness values of the corresponding output traces. For simplicity of notation, for a given sampled parameter vector  $x$ , we write  $\rho(x)$  to denote the robustness value of the corresponding output trace. And for a set  $S$  of sampled parameter vectors,  $\rho(S) = \min\{\rho(x) \mid x \in S\}$ . A parameter vector  $x$  is called a *falsifier* if  $\rho(x) < 0$ .

A rudimentary approach to search for a falsifier is repeatedly select randomly an unoccupied cell with uniform probability distribution and evaluate one point inside it. This way, we ensure that the cell occupancy always keeps increasing until we eventually find a robust bug, if it exists (see Lemma 1); however, this approach may not be efficient because the uniform search does not differentiate regions that are more likely to contain an input that falsifies from those that are less likely. Therefore, we propose to enhance it using two concepts:

1. *Using classification to bias random search.* We use robustness based classification to classify less falsifiable regions from more falsifiable regions. Then, the probability distribution of random samples is biased according to the coverage and robustness information in different regions.
2. *Combining global search and local search.* Local search approaches (such as Hill climbing, Gradient methods, Simulated annealing, Genetic algorithms) (see for example [24, 19]) can be very efficient if the search procedures are appropriately initialized. Finding good initializations constitutes a major difficulty that limits the efficiency of these approaches. In our framework, the classification based global

search provides useful hints at appropriate initializations for the local search. Indeed, the least robust points in regions with high potential of falsification can be used to initialize a number of local searches.

Thus, our falsification algorithm involves two phases. The first phase is a global search guided by hyperplane classifiers, coverage and the robustness information. Next is a local search phase which runs a number of local searches initiated at the least robust points of different regions formed by the classification process during the global search. We now explain the aforementioned ingredients of the algorithm.

#### 4.1 Classification using hyperplane subdivision

In the following, we say that two regions are *separate* if their intersection can only lie on their boundaries. Our classification problem can be intuitively described as follows: given a rectangle  $R$  representing a search space and a set  $S$  of sampled points in  $R$ , iteratively subdivide it, according to the robustness values of the sampled points, to obtain a rectangular partition, the elements of which have different average robustness levels.

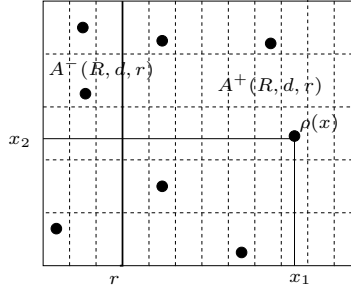
We define the average robustness of the set of samples  $S$  in  $R$  as  $\mu = \frac{\sum_{x \in S} \rho(x)}{|S|}$ . Our

objective is thus to separate a region of  $R$  having higher potential of containing low robustness samples. To this end, we define a hyperplane, in view of separating samples below the average robustness  $\mu$  from those above  $\mu$ . Obviously, such a hyperplane does not always exist, and we therefore choose a hyperplane that does this separation as best as possible. To address this problem, we draw inspiration from *soft margin support vector machines* [5], where hyperplanes are determined so that the misclassification error is minimized. In general, a misclassification error is defined according to the locations of the misclassified samples; however, in our approach we define a misclassification error that gives weight to the robustness values of the misclassified samples in addition to their locations. Furthermore, since it is easy to sample uniformly in rectangles, we only use axis-aligned hyperplanes, which generate only rectangular subregions. Otherwise, when allowing non-axis aligned classifiers, we generate polyhedral regions in which uniform random sampling as well as partition manipulation could be more expensive.

To explain the essence of our classification method, let us consider one rectangle  $R$  in the partition as the product of intervals  $R = [a_1, b_1] \times \dots \times [a_n, b_n]$ . Let  $S$  be the set of samples in  $R$ . We denote an axis-aligned hyperplane inside  $R$  by a tuple  $(d, r)$  where  $d \in \{1, \dots, n\}$  is the axis normal to the separating hyperplane, while  $r \in [a_d, b_d]$  is a coordinate at which the hyperplane is drawn. The hyperplane  $(d, r)$  subdivides  $R$  into two subrectangles  $A^-(R, d, r)$  and  $A^+(R, d, r)$  such that  $A^-(R, d, r) = [a'_1, b_1] \times \dots \times [a'_n, b_n]$  where  $a'_j = r$  if  $j = d$ , and  $a'_j = a_j$  otherwise; and  $A^+(R, d, r) = [a_1, b'_1] \times \dots \times [a_n, b'_n]$  where  $b'_j = r$  if  $j = d$ , and  $b'_j = b_j$  otherwise.

An ideal separation of samples below the average robustness from those above the average robustness by the hyperplane can be described in one of the following scenarios, identified by the following notion of *polarity*. Hyperplane  $(d, r)$  has polarity  $p = 1$  w.r.t  $S$ , if the left subrectangle  $A^-(R, d, r)$  contains all samples below the average  $\mu$ , while the right subrectangle  $A^+(R, d, r)$  contains all samples above the average  $\mu$ . Similarly,  $(d, r)$  has polarity  $p = -1$  w.r.t  $S$ , if  $A^+(R, d, r)$  contains all samples below  $\mu$ , while





**Fig. 1.** Classification by subdivision. The samples in a 2-dimensional rectangle  $R$  are represented by black points labeled with their robustness values;  $R$  is divided by the hyperplane  $(d, r)$  where the axis is  $d = 1$ , to minimise the misclassification error. This division produces two subrectangles  $A^-(R, d, r)$  (on the left) and  $A^+(R, d, r)$  (on the right).

$A^-(R, d, r)$  contains all samples above  $\mu$ . When an ideal separation as above is not feasible, we identify misclassified samples as follows.

**Definition 1.** A point  $x \in R$  is misclassified w.r.t. a hyperplane  $(d, r)$ , polarity  $p \in \{-1, 1\}$  and the sampled set  $S$ , if  $\text{sgn}(p(\rho(x) - \mu)(x_d - r)) = 1$  where  $x_d$  is the  $d^{\text{th}}$  coordinate of  $x$ , and  $\text{sgn}$  denotes the sign function.

For a misclassified sample, the misclassification error is measured according to its location and robustness value. If a misclassified sample is farther from the hyperplane, it is considered to entail a higher misclassification error. Also, since the classification is based on the average robustness, samples with robustness values farther from the average get higher weight in measuring the misclassification error. Accordingly, we define the misclassification error for a point  $x \in R$  w.r.t a hyperplane  $(d, r)$ , polarity  $p \in \{-1, 1\}$ , and a set of samples  $S$  as  $e_{d,r}(x, R, S, p) = \max\{p(\rho(x) - \mu)(x_d - r), 0\}$ . Then the total misclassification error is the sum of the misclassification errors of all the samples:

$$\Gamma_{d,r}(R, S, p) = \sum_{x \in S} e_{d,r}(x, R, S, p). \quad (4)$$

An appropriate hyperplane  $(d_*, r_*)$  traversing the rectangle  $R$ , chosen for the desired separation of  $S$ , is one that minimises the total misclassification error for either a positive or negative polarity, *i.e.*,

$$(d_*^{R,S}, r_*^{R,S}) = \arg \min_{r \in [a_d, b_d], d \in [n]} (\min\{\Gamma_{d,r}(R, S, -1), \Gamma_{d,r}(R, S, 1)\}). \quad (5)$$

We denote  $A_*^-(R, S) = A^-(R, d_*^{R,S}, r_*^{R,S})$  and  $A_*^+(R, S) = A^+(R, d_*^{R,S}, r_*^{R,S})$  as the subrectangles formed by dividing  $R$  by the above optimal hyperplane.

It is important to remark that in order for the classification to reflect the robustness distribution over the whole dense space, the number of samples should be sufficiently large. Henceforth, only rectangles in which the number of samples is not smaller than a (user-defined) threshold number, are subdivided as above. The classification procedure

takes as input a partition encoded as a list of  $k$  rectangles and the set of points in each respective rectangle. For each rectangle if the number of points is not smaller than the threshold  $K_c$ , the rectangle is subdivided by a hyperplane that minimizes the classification error. The rectangle is replaced by the left subrectangle and the right subrectangle is added to the list of rectangles. After all rectangles are considered for subdivision, the samples inside them are updated.

## 4.2 Global search

Each iteration of the global search performs 3 successive procedures:

1. Classification using hyperplanes, the goal of which is to partition the state space into regions with different robustness levels.
2. Coverage and robustness guided sampling of input signal parameters.
3. Singularity based sampling of input signal parameters inside rectangles containing very low robust samples. Here, we use *singularity* to refer to a partition element that contains a point in a low robustness range with low frequency of occurrence.

Note that the term *sampling* in the description of our method refers to the consecutive execution of three steps (1) defining the input signals from the sampled parameters, (2) simulating the system under the defined input signals, and (3) evaluating the robustness of the corresponding simulated output traces.

**Coverage and robustness based sampling** We randomly select a number of unoccupied cells, such that the probability of picking a cell in each rectangle is based on two components: coverage based probability and robustness based probability. Then the probability of cell sampling is determined as a weighted sum of the former components. Once a cell is sampled, a point is selected by a uniform sampling inside the cell.

*Coverage based probability distribution.* Let  $\{R_1, \dots, R_k\}$  be the set of rectangles of a partition of the parameter space. We now consider the collection of grid cells intersecting with  $R_i$ , that is  $\{\omega_j : \omega_j \cap R_i \neq \emptyset\}$ , and we index them as  $\beta^i = \{\beta_1^i, \dots, \beta_{q_i}^i\}$  where  $q_i$  is the number of such cells. Let  $D(R_i, S_i)$  be the vector denoting the distribution of samples  $S_i$  in cells of  $\beta^i$ , that is  $\forall j \in \{1, \dots, l\}$  ( $l$  is the total number of grid cells), the  $j^{th}$  component  $D_j(R_i, S_i) = |S_i \cap \beta_j^i|$ . Then the coverage based sampling probability in  $R_i$  is proportional to the number of unoccupied cells in this rectangle:

$$P_c^i = \frac{1 - H_c(D(R_i, S_i))}{\sum_{j=1}^m (1 - H_c(D(R_j, S_j)))}, \quad (6)$$

where  $H_c(D(R_i, S_i))$  is the *local cell occupancy* of  $R_i$ , i.e.,  $H_c(D(R_i, S_i)) = \frac{\log(N_{c_i})}{\log(l_i)}$ . where  $l_i$  is the number of grid cells intersecting with  $R_i$  and  $N_{c_i}$  is the number of unoccupied cells intersecting with  $R_i$ .

**Robustness based probability distribution.** A probability distribution takes into consideration the average robustness as well as the potential reduction in robustness below the average. The potential reduction in robustness below the average is defined as  $\lambda_i = \frac{1}{|S_i|} \sum_{x \in S_i} \max(\mu_i - \rho(x), 0)$ . Then a potentially reduced robustness value below the average is  $\theta_i = \mu_i - \lambda_i$ . Then we define a robustness based probability in a rectangle  $R_i$  as inversely proportional to  $\theta_i$ , as follows.

$$P_r^i = \frac{\frac{1}{\theta_i}}{\sum_{j=1}^m \frac{1}{\theta_j}}. \quad (7)$$

**Sampling probability distribution.** The probability distribution for sampling is a weighted sum of the probability based on robustness and the probability based on coverage. The weight given to either probability is a user defined constant. Let the weight assigned to the robustness based probability be denoted by  $w_r$  such that  $w_r \in [0, 1]$ . Then, the overall probability of sampling in rectangle  $R_i$  is

$$P_t^i = w_r P_r^i + (1 - w_r) P_c^i. \quad (8)$$

**Singularity based sampling** Certain rectangles may contain samples whose robustness is very low compared to the lowest robustness values in other rectangles. We refer to them as *singular samples*, which we heuristically define as follows.

Let  $\gamma = \{\gamma_1, \dots, \gamma_k\}$  be the vector of lowest robustness values in each rectangle, defined as  $\gamma_i = \min_{x \in S_i} \rho(x)$ . The mean of  $\gamma$  and the average deviation below the mean are respectively defined as  $\mu_\gamma = \frac{\sum_{i=1}^k \gamma_i}{k}$  and  $\lambda_\gamma = \frac{\sum_{i=1}^k \max(0, (\mu_\gamma - \gamma_i))}{k}$ . If the robustness of a sample is less than  $\lambda$  then it is an indication that the sample may be close to a counterexample. Also, samples with very low frequency and sufficiently low robustness are also considered singular. To select such rare samples, we use the following heuristic. If  $\gamma$  were a large set of random samples selected from a normal distribution, then less than 15% of the samples tend to lie below the value  $\mu_\gamma - 3\lambda_\gamma$ . Although the actual set of samples in  $\gamma$  may not follow the pattern of a normal distribution, this also can be used as a heuristic to define a singular sample.

**Definition 2.** A point  $x \in \bigcup_{i=1}^k S_i$  for which  $\rho(x) \leq \max(\mu_\gamma - 3\lambda_\gamma, \lambda)$  is called a *singular sample*.

We call the rectangles containing singular samples as *singular rectangles*. Since the frequency of singular samples can be very small, the robustness based probability in (7) may not give adequate weight to singular rectangles. So, we have to perform additional sampling in the singular rectangles.

**Overall global search** Suppose that we have a partition of rectangles  $R_1, \dots, R_k$  containing sets of samples  $S_1, \dots, S_k$ , respectively. Let  $C_i$  be the set of unoccupied cells intersecting with a rectangle  $R_i$ , i.e.,  $C_i = \{\omega_j \in \omega : \omega_j \cap R_i \neq \emptyset \wedge \omega_j \cap S_i = \emptyset\}$ . Let  $N$  be the number of samples to be added during probabilistically biased random

sampling. We compute the probability distribution of sampling among different rectangles  $P^t$ , where the user defines a weight  $w_r$  given to the robustness based probability distribution  $P^r$ . Then we select  $(\min\{\max\{1, \lfloor P_i^t N \rfloor\}, |C_i|\})$  number of cells among  $C_i$  and sample one point in each cell. Note that if there is an unoccupied cell in a rectangle, then at least one sample is added to each rectangle irrespective of the probability  $P_i^t$ . Then update the sets of samples  $S_1, \dots, S_k$ , by adding the new samples and also the sets of unoccupied cells  $C_i \forall i \in \{1 \dots k\}$ .

Next, we perform sampling in each of the singular rectangles as follows. Let  $R_j$  be a singular rectangle, currently containing the set  $S_j$  of samples and a set  $C_j$  of unoccupied cells. Then we select  $\min\{\max\{K_c - |S_j|, 0\}, |C_j|\}$  cells among the unoccupied cells  $C_i$ . Therefore, in the next iteration  $R_j$  contains at least  $K_c$  samples (if it has unoccupied cells) and is consequently subdivided. The procedure is repeated in each iteration until the time limit  $\bar{T}_g$  on global search is reached. Alternatively, we can also set a limit on the total number of samples for which robustness is evaluated. If we have not falsified yet, then we perform a number of local searches initialized at the lowest robustness samples of all the separate subrectangles, as described below.

### 4.3 Local search

Suppose that we have  $K$  subrectangles  $R_1, \dots, R_K$  after running global search for  $\bar{T}_g$  time. Let  $L$  be the set of the lowest robustness points of different rectangles. If the property is not yet falsified, then we use the lowest robustness points of different rectangles to initialize a local search based falsification algorithm. In our implementation, we used the state-of-the-art *Covariance matrix adaptive evolutionary search* (CMA-ES) algorithm [21] in the local search phase. The essence of the CMA-ES algorithm can be briefly described as follows. It is a randomized black box method which selects samples based on a multivariate normal distribution having a mean and a covariance matrix as parameters. Based on the robustness of a population of points evaluated in an iteration, the mean and covariance matrix of the search distribution are updated for the next sampling iteration. The procedure generally converges to a locally optimum point or finds a counterexample. It may happen that the set of sampled points do not contain enough information to derive a reliable estimation of a covariance matrix for an efficient update. Therefore, good initializations of the mean and covariance matrix are crucial. In our algorithm, the global search provides initialization guidance as follows. We have a number of subrectangles formed by classification, that contain sets of samples. So, we can initialize in one the following ways: (1) Each of the lowest robustness points of different rectangles, i.e., the points in  $L$  can be selected for initialization with covariance as identity (the order of selection is according to their robustness values, with the lowest robustness tested first); (2) The mean and covariance are initialized as that of those points in  $L$  whose robustness is less than the average robustness of samples in  $L$ . (3) The mean and covariance are initialized as that of all the points in  $L$ . With such initialization guidance from our global search procedure described earlier, this local search procedure can become more effective in falsification.

#### 4.4 Overall falsification algorithm

The overall falsification algorithm consists of iteratively doing global search for a threshold time and then doing local search. We give an outline of the algorithm below.

- Step 1: *Initialization*. In the first step, we evaluate the robustness of  $N$  uniformly selected points in the search space  $R$  and store them as a set of samples  $S$ .
- Step 2: *Global search phase*. We perform a number of global search iterations until a time limit is attained. Each global search iteration consists of the following three sequential procedures: (i) The first step is classification, where new rectangles are constructed by classifying and consequently subdividing the existing rectangles that contain more than a threshold  $K_c$  of samples. (ii) The second step involves probabilistically biased sampling based on the coverage and robustness values of samples in different rectangles. (iii) The third step is singularity based sampling. The detailed methods in the second and third procedures are explained earlier in Section 4.2.
- Step 3: *Local search phase*. If no counterexample is found during Step 2, then we perform the local search based on the set of low robustness points in different rectangles. The specific procedure is explained earlier in Section 4.3.
- Step 4: If not falsified during Step 3, then go to Step 2 to continue global search iterations.

We can now state an important completeness property of our overall falsification algorithm for the class of the systems (2) satisfying the assumption that the mapping  $\hat{\Phi}$  is continuous in the nominal parameters and the input signal parameters.

**Theorem 1.** *If an  $\epsilon$ -robust counterexample exists, then there exists a grid cell size  $\Delta$  and a global search time  $\bar{T}_g$  so that our algorithm finds a counterexample.*

*Sketch of proof.* The theorem can be directly established from Lemma 1. Indeed, the condition in this lemma is always satisfied by our algorithm since, by construction, after each iteration the cell occupancy of the samples always strictly increases. So, for sufficient  $\bar{T}_g$ , the falsification is guaranteed if an  $\epsilon$ -robust counterexample exists.

## 5 Experimental Results

In our experiments, we compare the performance of a MATLAB implementation<sup>5</sup> with the following standard approaches: CMA-ES, Simulated Annealing, Global Nelder-Mead algorithm implementations (integrated in Breach [9]), and the S-TaLiRo tool [2] by setting Simulated Annealing as optimization algorithm<sup>6</sup>. experiments were performed on a computer with 1.4GHz processor with 4GB RAM, running MATLAB R2015 64-bit version. Also, we compare with a random sampling method, where in each iteration

<sup>5</sup> We use the robustness evaluation function from the Breach toolbox available in October 2016, on the site [https://people.eecs.berkeley.edu/~donze/breach\\_page.html](https://people.eecs.berkeley.edu/~donze/breach_page.html)

<sup>6</sup> We used the latest version available in October 2016, on the site <https://sites.google.com/a/asu.edu/s-taliro/home>.

a pseudo-randomly selected point is tested only if it falls in a grid cell wherein no other point has been previously tested. The grid used in this method is the same as the grid chosen in our falsification approach. We will call this method as grid based random sampling, for the sake of reference during comparison.

## 5.1 Automotive Powertrain Control

We consider a Simulink model of a closed loop of an Automotive Powertrain Control subsystem (PTC). The model contains a representation of an internal combustion engine and an embedded software controller for the air-to-fuel ratio within the engine (see [11] for more details). Here, we focus on the input-output behavior, considering the internal model as a blackbox. The model has three input signals, Pedal Angle Engine Speed and Sensor Offset. The air-to-fuel (A/F) ratio, denoted by  $\eta$ , is an output signal for which the following safety requirement was stated in [11]:  $\phi = \square_{[5,10]} (\eta < 0.5)$ .

**Input signal settings.** Compared to [11], we consider a smaller input range for the Pedal Angle as  $[0, 40]$  and fix the Engine Speed and Sensor Offset as 1000 and 1, respectively. Reducing the ranges makes the properties more robust and consequently difficult to falsify. The time horizon is 50s. We use piecewise constant signal for testing, where the Pedal Angle is parameterized by 10 uniformly spaced control points in the time horizon. Thus, we have a 10 dimensional search space  $X$ .

**Algorithm setting.** For our algorithm, the threshold number of samples for hyperplane classification  $K_c$  is 100. The global search time is  $\bar{T}_g = 2000$  seconds. The local search is initialized with the lowest robustness point found during global search and allowed to run until falsification. Cell partitioning  $\omega$  consists of hypercubes of side length  $\epsilon = 4$ . We consider equal weightage for robustness based probability and coverage based probability for sampling during global search, i.e.,  $w_r = 0.5$ .

**Results.** Our algorithm (classification guided global search + local search) successfully found a counterexample in less than 3000 seconds for all seeds. As an estimate of the classification frequency, the final number of separate rectangles constructed for while testing the first seed were 30. In comparison, the tool S-TaLiRo could falsify but took 4481 seconds. The grid based random sampling found a falsifier for only the seeds 15000 and 20000, but failed to do so on the other seeds before maximum time limit was reached. The other methods were not successful in finding a falsifier within the default stopping time of 5000 seconds. Both the CMA-ES and Nelder-Mead became stuck without reduction in robustness value until the default stopping time was reached. The results are presented in Table 1. We note that for any fixed seed for random sampling, these results are reproducible.

## 5.2 Automatic transmission

We consider the benchmark model of an Automatic Transmission control system, which appeared in [20]<sup>7</sup>. The system has two input signals, called throttle and break, respectively, and two output signals, called the engine speed, denoted  $w$  (RPM), and the vehicle speed, denoted  $v$  (mph). The property states that if the engine speed stays below a

<sup>7</sup> The model and property description of this benchmark is available at the site of the workshop Applied Verification for Continuous and Hybrid Systems, ARCH 2014-2015, <http://cps-vo.org/node/12116>

value  $\bar{w}$ , then the vehicle speed  $v$  does not exceed a threshold  $\bar{v}$  within 10 seconds. We specify the values of  $\bar{w}$  and  $\bar{v}$  to be 2520 and 50, respectively, which gives the following STL property:  $\phi = \neg ((\Diamond_{[0,10]} v > 50) \wedge (\Box w \leq 2520))$  [20].

**Input signal and parameter settings.** Initially, the vehicle is at rest, when  $v = 0$  and  $w = 0$ . For the input signals, we consider smaller ranges than specified in [20], which makes the property  $\phi$  more robust. Henceforth, the throttle signal is allowed to vary between  $[35, 100]$  and the brake is allowed to vary between  $[0, 40]$ . The time horizon is set to 30 seconds. We use piecewise constant input signals for testing, where the throttle signal is parametrized by 7 control points and the brake has 3 control points. Thus, we have a 10 dimensional search space.

**Algorithm setting.** The threshold number of samples of hyperplane classification  $K_c$  is 70. The global search time is  $\bar{T}_g = 500$  seconds, while maximum time for local search is  $\bar{T}_l = 2000$  seconds. Cell partitioning  $\omega$  consists of hypercubes of side length  $\epsilon = 4$ . We consider equal weightage for robustness based probability and coverage based probability for sampling during global search, i.e.,  $w_r = 0.5$ .

**Results.** Our algorithm (classification guided global search with local search) successfully found a counterexample in less than 2000 seconds for all randomly chosen seeds. As an indication of the number of classification operations that occurred, the final number of separate rectangles constructed for while testing the first seed were 31. In comparison, the CMA-ES found a falsifier for two seeds 5000 and 15000 within 2000 seconds but failed to do so on the other seeds. The other methods were not successful in finding a falsifier within the default stopping time of 3000 seconds. For this example, S-TaLiRo became stuck around a local optimum without any significant reduction in robustness value. The results are presented in Table 1. We note that for any fixed seed for random sampling, these results are reproducible.

### 5.3 Industrial example

We present results for an air path controller for an automotive fuel cell (FC) application. The system contains an FC stack that generates electrical power to provide torque to the vehicle drivetrain. The system is composed of an air compressor and the air path through the FC stack. The system takes as input requested current from the stack and ambient temperature. The outputs are desired air flow rate and the measured air flow rate through the FC stack. The goal is for the stack air flow rate to maintain accurate regulation when current request “disturbances” are presented to the system. System performance (called *responsiveness*) crucially depends on accurate and timely regulation of the air flow to the commanded reference. The corresponding specification for the system can be described informally as follows: when there is a step input of current request, there is a rise-time requirement on the output air flow that should be satisfied. Details about the system and the specifications are proprietary and so are suppressed here.

We analyze a Simulink model of the FC system, which contains representations of the FC system along with its controller. The model is complex, containing several thousands of Simulink blocks; simulations over the selected time horizon are expensive to perform, each taking approximately 1 to 2 minutes. The MATLAB implementation of the hyperplane classification algorithm with local search is applied to the model, and the results are compared to the same algorithms used in Sections 5.1 and 5.2.

**Table 1.** Experimental results

Solver	Seed	Computation time (secs)		Falsification	
		PTC	Aut. Trans	PTC	Aut. Trans.
Hyperplane classification + CMA-ES-Breach	0	2891	996	✓	✓
	5000	2364	1382	✓	✓
	10000	2101	1720	✓	✓
	15000	2271	1355	✓	✓
CMA-ES-Breach	0	T.O. (5000)	T.O. (2000)		
	5000	T.O. (5000)	1302		✓
	10000	T.O. (5000)	T.O. (2000)		
	15000	T.O. (5000)	1325		✓
Grid based random sampling	0	T.O. (5000)	T.O. (2000)		
	5000	T.O. (5000)	T.O. (2000)		
	10000	3766	T.O. (2000)	✓	
	15000	268	T.O. (2000)	✓	
S-TaLiRo (Simulated Annealing)		4481	T.O. (3000)	✓	
S-TaLiRo (Simulated Annealing)		4481	default stopping (3300)	✓	

*T.O.*: Exceeded indicated time out limit.

*Seed*: Index for a sequence of random numbers in MATLAB. *Solver*: Algorithm used for falsification. *Computation time*: Amount of time (in seconds) until falsification or default stopping after the time limit in parentheses. Computation time is reported for a computer with 1.4GHz processor and 4GB RAM, running MATLAB R2015 64-bit version. *Falsification*: Boolean variable indicating whether the algorithm could falsify the property.

For our method, we performed the tests using two different cell partitions. Cell partition A is large and corresponds to a small number of grid elements; cell partition B is smaller (each dimension of the search space is 1/5 the size of the grid elements in partition A). Thus, partition B corresponds to a significantly larger number of grid elements.

Table 2 provides the results. As can be seen in the table, using cell partition A with our method performs much better than with partition B. This can be attributed to the fact that, for partition A, the classification phase of the search spends less time in regions close to regions that have already been explored, as compared to partition B. This demonstrates that the selected cell partition size has a significant impact on the performance of our technique.

Also, the table shows that the CMA-ES fails to find falsifying behaviors in 2 of the 4 cases, which demonstrates better performance than our technique using partition B but poorer performance than our technique using partition A. The uniform random sampling approach is able to find falsifying traces in all but one case, and the computation times for the successful cases are comparable to our technique using partition A, though we note that the computation times for our technique are lower than the uniform ran-



dom method, for the cases where falsifying traces are found. The S-TaLiRo approach fails to find falsifying traces in 2 of the 4 cases, which is less than the number of times our technique is successful, using partition A. The Nelder-Mead algorithm is able to identify a falsifying trace in about 25 minutes, which is longer than the 3 successful cases of our technique, using partition A.

The above results show mixed results for our technique for this example, as compared to the other falsification approaches. This could be due to any of several factors. We observe that for this example, comparing against the falsification techniques that we selected, only a relatively small number of simulations are required to find falsifying traces, when they are found at all. This may suggest that either the model is not robust, in the sense that there may be many disconnected regions in the search space that correspond to falsifying behaviors, or that the robustness function is rather monotone or simple. It may be that for systems with these qualities, the benefits provided by the hyperplane classification approach are outweighed (or at least offset) by the overhead that it requires.

## 6 Conclusions

We have presented a novel falsification algorithm that maintains a balance between convergence towards low robustness points and enhancing global coverage. We accomplish this by intelligently subdividing the search space and subsequently biasing the density of random sampling in different sub-regions. For the subdivision, we use hyperplane classifiers akin to support vector machines, which tries to focus effort on low robustness regions of the search space. We demonstrated the efficiency of our algorithm by falsifying properties on benchmark examples, which other approaches failed to falsify. Also, we demonstrated that the approach could be applied to industrial systems by describing a successful application on an automotive hydrogen fuel cell example. Future work includes investigating new coverage measures, such as the combinatorial entropy notion from the domain of physics to measure the degree of randomness in the distribution of points. In addition, global search and local search can be done in a multi-resolution manner, that is if local search leads to a promising region, global search can then be done within the region using a more refined grid.

## References

1. M. Althoff and B. Krogh. Zonotope bundles for the efficient computation of reachable sets. In *Decision and Control and European Control Conference (CDC-ECC), 2011 50th IEEE Conference on*, pages 6814–6821, Dec 2011.
2. Y. Annapureddy, C. Liu, G. E. Fainekos, and S. Sankaranarayanan. S-TaLiRo: A Tool for Temporal Logic Falsification for Hybrid Systems. In *TACAS*, pages 254–257, 2011.
3. O. Bouissou, E. Goubault, S. Putot, K. Tekkal, and F. Védreine. HybridFluctuat: A static analyzer of numerical programs within a continuous environment. In *Computer Aided Verification, 21st International Conference, CAV 2009, Grenoble, France, June 26 - July 2, 2009. Proceedings*, pages 620–626, 2009.
4. X. Chen, E. Abraham, and S. Sankaranarayanan. Flow\*: An Analyzer for Non-Linear Hybrid Systems. In *CAV*, 2013.

**Table 2.** Results for Fuel Cell Example.

Solver	Seed	Computation time (sec.)	Falsification
Hyperplane classification + CMA-ES-Breach (Cell partition: A) <sup>†</sup>	1	406	✓
	2	1383	✓
	3	T.O.	
	4	794	✓
Hyperplane classification + CMA-ES-Breach (Cell partition: B) <sup>†</sup>	1	409	✓
	2	T.O.	
	3	T.O.	
	4	T.O.	
CMA-ES Breach <sup>‡</sup>	1	314	✓
	2	1418	
	3	T.O.	
	4	1316	✓
Uniform random <sup>‡</sup> sampling	1	396	✓
	2	786	✓
	3	2241	✓
	4	T.O.	
S-TaLiRo (Simulated Annealing) <sup>‡</sup> sampling	1	310	✓
	2	T.O.	
	3	671	✓
	4	T.O.	
Global Nelder-Mead-Breach <sup>‡</sup>		1501	✓

T.O.: Exceeded time out limit of 2700 seconds.

<sup>†</sup>: Times reported are from machines running Dell Precision, with a Xeon processor (2.13GHz), with 24GB of RAM, running a 64 bit version of Windows 7 Ultimate, SP1.

<sup>‡</sup>: Times reported are from machines running Dell Precision, with a Xeon processor (2.3GHz), with 64GB of RAM, running a 64 bit version of Windows 7 Ultimate, SP1.

5. C. Cortes and V. Vapnik. Support-vector networks. *Mach. Learn.*, 20(3):273–297, Sept. 1995.
6. T. Dang and T. Nahhal. Coverage-guided test generation for continuous and hybrid systems. *Formal Methods in System Design*, 34(2):183–213, 2009.
7. J. V. Deshmukh, X. Jin, J. Kapinski, and O. Maler. Stochastic local search for falsification of hybrid systems. In *Automated Technology for Verification and Analysis - 13th International Symposium, ATVA 2015, Shanghai, China, October 12-15, 2015, Proceedings*, pages 500–517, 2015.
8. T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1):31–71, 1997.
9. A. Donzé. Breach, a toolbox for verification and parameter synthesis of hybrid systems. In *CAV*, pages 167–170, 2010.
10. A. Donzé and O. Maler. Robust satisfaction of temporal logic over real-valued signals. In *FORMATS*, pages 92–106, 2010.
11. T. Dreossi, T. Dang, A. Donzé, J. P. Kapinski, X. Jin, and J. V. Deshmukh. Efficient guiding strategies for testing of temporal properties of hybrid systems. In *NASA Formal Methods - 7th*

- International Symposium, NFM 2015, Pasadena, CA, USA, April 27-29, 2015, Proceedings*, pages 127–142, 2015.
12. T. Dreossi, T. Dang, and C. Piazza. Parallelotope bundles for polynomial reachability. In *Proceedings of the 19th International Conference on Hybrid Systems: Computation and Control, HSCC 2016, Vienna, Austria, April 12-14, 2016*, pages 297–306, 2016.
  13. J. M. Esposito, J. Kim, and V. Kumar. Adaptive rrts for validating hybrid robotic control systems. In *Algorithmic Foundations of Robotics VI*, pages 107–121. Springer Berlin Heidelberg, 2005.
  14. G. Fainekos and G. J. Pappas. Robustness of temporal logic specifications. In *Proceedings of FATES/RV*, volume 4262 of *LNCS*, pages 178–192. Springer, 2006.
  15. C. Fan, B. Qi, S. Mitra, M. Viswanathan, and P. S. Duggirala. Automatic reachability analysis for nonlinear hybrid models with C2E2. In *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part I*, pages 531–538, 2016.
  16. G. Frehse, C. Le Guernic, A. Donzé, S. Cotton, R. Ray, O. Lebeltel, R. Ripado, A. Girard, T. Dang, and O. Maler. SpaceX: Scalable verification of hybrid systems. In *CAV*, pages 379–395, 2011.
  17. G. M. Fung, O. L. Mangasarian, and J. W. Shavlik. Knowledge-based support vector machine classifiers. *Advances in neural information processing systems*, pages 521–528, 2002.
  18. S. Gao, J. Avigad, and E. M. Clarke.  $\delta$ -complete decision procedures for satisfiability over the reals. In *J. Automated Reasoning*, pages 286–300, 2012.
  19. H. Hoos and T. Sttze. *Stochastic Local Search: Foundations & Applications*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2004.
  20. B. Hoxha, H. Abbas, and G. E. Fainekos. Benchmarks for temporal logic requirements for automotive systems. In *1st and 2nd International Workshop on Applied verification for Continuous and Hybrid Systems, ARCH@CPSWeek 2014, Berlin, Germany, April 14, 2014 / ARCH@CPSWeek 2015, Seattle, WA, USA, April 13, 2015.*, pages 25–30, 2014.
  21. C. Igel, T. Sutton, and N. Hansen. A computational efficient covariance matrix update and a (1+1)-CMA for evolution strategies. In *Proceedings of the 8th annual conference on genetic and evolutionary computation GECCO*, pages 453–460. ACM, 2006.
  22. J. Kuřátko and S. Ratschan. Combined global and local search for the falsification of hybrid systems. In *International Conference on Formal Modeling and Analysis of Timed Systems*, pages 146–160. Springer, 2014.
  23. O. Maler and D. Nickovic. Monitoring temporal properties of continuous signals. In *FORMATS/FTRTFT*, pages 152–166, 2004.
  24. S. J. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 2 edition, 2003.
  25. P. Skrch. A coverage metric to evaluate tests for continuous-time dynamic systems. *Central European Journal of Engineering*, 1(2):174180, 2011.
  26. R. Testylier and T. Dang. NLTOOLBOX: A library for reachability computation of nonlinear dynamical systems. In *Automated Technology for Verification and Analysis - 11th International Symposium, ATVA 2013, Hanoi, Vietnam, October 15-18, 2013. Proceedings*, pages 469–473, 2013.