# Reflections on Surrogate-Assisted Search-Based Testing: A Taxonomy and Two Replication Studies based on Industrial ADAS and Simulink Models

Shiva Nejati[a,*], Lev Sorokin[b], Damir Safin[b], Federico Formica[c],
Mohammad Mahdi Mahboob[c], Claudio Menghi[c,d]

[a]*University of Ottawa, Canada, Ottawa, K1N6N5, Canada*
[b]*Fortiss, Munich, Germany*
[c]*McMaster University, 1280 Main St W, Hamilton, ON L8S4L8, Canada*
[d] *University of Bergamo, via Salvecchio 19, Bergamo, 24129, Italy*

**Abstract**

*Context.* Surrogate-assisted search-based testing (SA-SBT) aims to reduce the computational time for testing compute-intensive systems. Surrogates enhance testing techniques by improving test case generation focusing the testing budget on the most critical portions of the input domain. In addition, they can serve as approximations of the system under test (SUT) to predict test results instead of executing the tests on compute-intensive SUTs.

*Objective.* This article reflects on the existing SA-SBT techniques, particularly those applied to system-level testing and often facilitated using simulators or complex test beds. Recognizing the diversity of heuristic algorithms and evaluation methods employed in existing SA-SBT techniques, our objective is to synthesize these differences and present a comprehensive view of SA-SBT solutions. In addition, by critically reviewing our previous work on SA-SBT, we aim to identify the limitations in our proposed algorithms and evaluation methods and to propose potential improvements.

*Method.* We present a taxonomy that categorizes and contrasts existing SA-SBT solutions and highlights key research gaps. To identify the evalua-

---

*Corresponding author

*Email addresses:* snejati@uottawa.ca (Shiva Nejati), sorokin@fortiss.org (Lev Sorokin), safin@fortiss.org (Damir Safin), formicaf@mcmaster.ca (Federico Formica), mahbom2@mcmaster.ca (Mohammad Mahdi Mahboob), claudio.menghi@unibg.it (Claudio Menghi)

tion challenges, we conduct two replication studies of our past SA-SBT solutions: One study uses industrial advanced driver assistance system (ADAS) and the other relies on a Simulink model benchmark. We compare our results with those of the original studies and identify the difficulties in evaluating SA-SBT techniques, including the impact of different contextual factors on results generalization and the validity of our evaluation metrics.

*Results.* Based on our taxonomy and replication studies, we propose future research directions, including re-considerations in the current evaluation metrics used for SA-SBT solutions, utilizing surrogates for fault localization and repair in addition to testing, and creating frameworks for large-scale experiments by applying SA-SBT to multiple SUTs and simulators.

## 1. Introduction

To test complex and emerging systems at scale, the use of simulators is necessary. Simulators provide a flexible, effective, and efficient infrastructure to ensure the safety of complex systems compared to testing deployed systems in their operational environment (e.g., vehicles' on-road testing). However, executing each simulation scenario still takes a non-negligible amount of time, particularly when simulators need to be real-time such as in the case of simulators for autonomous vehicles. Since simulators' input spaces are very large and high-dimensional, we require techniques to make simulation-based testing scalable. For example, non-trivial simulations of an industrial model of a satellite system, capturing the satellite behavior for 24h, take, on average, around 84 minutes ($\sim$1.5 hours) Menghi et al. (2020). To test this satellite model for all of its parameters and inputs, we need to execute hundreds of such simulations, which can take months or even years to complete.

Search-based testing (SBT) has been traditionally used as an effective and efficient guidance for test generation performed at system-level Zeller (2017). SBT algorithms provide flexibility in representing the input search space and choosing a strategy to traverse a subset of the search space. To scale SBT to handle compute-intensive cyber-physical systems (CPS) with large and multi-dimensional input spaces, a promising approach is to combine SBT

with *surrogate models* Jin (2011) that enable effectively guiding the search exploration.

Surrogate models can enhance SBT in two ways: (1) They identify the most critical regions of the input domain (i.e., the regions that include tests revealing critical faults) and sample test inputs that are in these regions. (2) They predict simulation results instead of computing them. In both cases, surrogate models are generated and iteratively refined based on the intermediary test executions. As surrogate models improve through successive refinements, their predictions can help improve the search guidance, as described by the first use case above. Alternatively, their predictions are used in lieu of actual test executions, as suggested by the second use case above, to reduce the search computation time by executing fewer tests.

This article reflects on four papers published between 2014 and 2021 Matinnejad et al. (2014); Abdessalem et al. (2016, 2018a); Menghi et al. (2020) that propose surrogate-assisted search-based testing (SA-SBT) techniques for autonomous systems. The papers evaluated their proposed testing techniques on industrial Advanced Driving Assistance Systems (ADAS) Abdessalem et al. (2016, 2018a), a satellite system, and an open-source benchmark of Simulink models Menghi et al. (2020) as well as an industrial controller from the automotive domain Matinnejad et al. (2014). Since these papers are the main subjects of our analysis, we refer to these four papers as the *subject papers* hereafter in this article.

The subject papers apply SBT enhanced using surrogate models to system-level testing and demonstrate that this combination is effective for complex, emerging case studies such as self-driving and satellite systems. These subject papers complement existing SBT techniques that focus on testing software code and unit-level testing, and the research on testing and verification of CPS based on formal methods Clarke et al. (2001). These four papers, while studying CPS with different behaviours and varying degrees of autonomy and complexity, motivate the need for SA-SBT techniques as follows: In the context of CPS and given the safety-critical nature of these systems, system-level testing is typically performed using simulators. CPS simulators, particularly those that are real-time, are compute-intensive, and their input spaces are large and high-dimensional. SA-SBT techniques are then proposed to scale simulation-based testing for CPS.

The four subject papers adopt an empirical approach and evaluate their SA-SBT techniques using case studies from different domains. The techniques that these papers propose, although all can be categorized as SA-SBT,

vary in several details. For example, three of the subject papers Matinnejad et al. (2014); Abdessalem et al. (2016); Menghi et al. (2020) use surrogates to reduce the computation time of computing fitness functions, while the fourth paper Abdessalem et al. (2018a) uses surrogates to guide the search exploration more effectively. Given the differences in case studies and proposed techniques, it is uncertain how the results can be generalized to similar systems. In this reflection paper, we review our previous research and attempt to answer two questions: (1) What are the different types and variations of SA-SBT techniques? and (2) To what extent can we reproduce the results of these four subject papers? To answer the first question, we propose a taxonomy for SA-SBT techniques and use it to categorize the subject papers and some recent papers that build on our research. To answer the second question, we report two replication studies that attempt to reproduce the results of two of the subject papers Abdessalem et al. (2018a); Menghi et al. (2020). In particular, one replication study applies the SA-SBT solution proposed by Abdessalem et. al. Abdessalem et al. (2018a) to an industrial ADAS. For the second replication study, we present the results obtained by ARIsTEO Menghi et al. (2020), an approximation-refinement testing technique based on surrogate models, in three consecutive editions (2020 Ernst et al. (2020), 2021 G. Ernst et al. (2021), and 2022 Ernst et al. (2022)) of the ARCH competition ARCH. The ARCH competition is a friendly yearly competition between testing tools for continuous and hybrid systems ARCH.

We use Kolb's model of experiential learning Dybå et al. (2014) as a framework for our reflection study, which includes a taxonomy and two replication experiments. The Kolb model, shown in Figure 1, outlines the steps of experiential learning and emphasizes the role of reflection in turning direct experiences into abstract concepts and knowledge that drive further research. The process starts with obtaining concrete experience, followed by reflective observations in light of existing knowledge. Reflection should focus on any discrepancies between existing knowledge and the new experience. From reflections, the researcher develops observations into abstract concepts and generalizations, and tests these through active experimentation.

We structure our reflection study following the experiential learning loop in Figure 1. We develop a taxonomy (Section 2) that provides a set of abstract concepts and categories to help researchers classify and compare different SA-SBT research. We use the taxonomy to reflect on the current state of research and practice in the SA-SBT field, and to identify current research gaps (Section 3). We replicate two of our subject papers Abdessalem
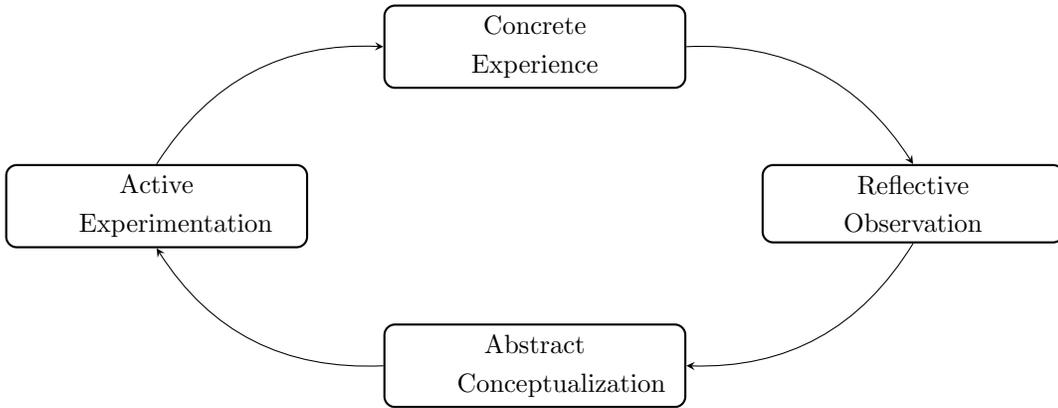
4

Figure 1: Experiential learning cycle.

et al. (2018a); Menghi et al. (2020) and present our observations based on the replication studies (Section 4). Our findings, which include the taxonomy, the categorization of the state-of-art along with the taxonomy dimensions, and two replication studies lead us to identify a number of shortcomings and challenges that impact the research on SA-SBT. We discuss future research directions that the software testing researchers can embark on to improve and unify the research in this field. Section 5 presents our recommendations for future research which are listed as follows: (1) standardizing metrics for evaluating SA-SBT solutions, (2) using surrogates to improve fault localization and repair, (3) creating frameworks for large-scale and consistent application of SA-SBT across multiple SUTs and simulators, and (4) developing benchmarks and organizing international competitions to promote the continuous replication of research tools.

## 2. Taxonomy

This section presents our taxonomy for SA-SBT research and classifies the existing research along these categories. Figure 2 presents a mind map diagram representing our taxonomy. The diagram visualizes the concepts captured by our taxonomy and their relations. SA-SBT (the center node of the mind map diagram) is connected to eight categories representing relevant problems to be considered in the design of SA-SBT techniques (purpose, usage, subject, design, type, context, scope, and evaluation metrics). Each

problem is connected with sub-categories nodes representing relevant solutions and parameters to be considered when addressing these problems. In the following, we first describe a brief background on search-based testing (SBT). We then describe the categories of our taxonomy and discuss the sub-categories related to each of these categories.



Figure 2: A Taxonomy for Surrogate-Assisted Search-Based Testing (SA-SBT).

**Background.** SBT relies on metaheuristic search Luke (2013) to generate a limited and effective set of test cases. Briefly, SBT algorithms work by building a set of initial candidate tests, iteratively applying tweak operations to one or more candidates to create new ones, and using a fitness function to decide whether to keep a candidate test for future iterations or to discard it.

**Purpose.** It refers to the way surrogates are combined with the metaheuristic search algorithm underlying an SBT technique and the search elements that they aim to improve. We identify three different purposes for surrogates: (1) *fitness computation*, (2) *test generation guidance*, and (3) *individuals selection*.

- *Fitness computation.* Surrogates reduce the fitness computation time by approximating the procedure for computing fitness values that often requires simulating the SUT.

- *Test generation guidance.* Surrogates replace or augment search operators (a.k.a. tweak operators). Specifically, they help identify the regions of the search space that include fitter individuals. New candidates can thus be generated from within these regions.

- *Individuals selection.* Decisions during the search process can be replaced with decisions guided by surrogates. For example, surrogates can be used to select fitter individuals for the initial population at the start.

**Usage.** It refers to the way surrogates improve the test generation approach of an SBT technique. We can use *individual* surrogates or an *ensemble* of surrogates. Surrogates may completely replace the SUT for test generation, or they may filter poor candidate test cases, while the more promising candidates are evaluated on SUT. We identify three different purposes for surrogates: (1) *test case generation using single surrogates*, (2) *test case filtering using single surrogates*, and (3) *test case generation using ensemble surrogates*.

- *Test case generation using single surrogates.* Surrogates are used to generate candidate test cases without any need to execute the original SUT during the test generation. To determine if the generated tests should be kept or discarded, they are assessed on the original SUT only at the end.

- *Test case filtering using single surrogates.* Surrogates are used to prune candidates that have no chance of surviving the search selection process. They only discard poor test cases. The remaining test cases are assessed on the original SUT.

- *Test case generation using ensemble surrogates.* An ensemble of surrogates is used to assess individual test cases. A voting mechanism is used to determine, for a given test case, if the surrogates are precise enough, i.e., the surrogates produce the same (or close) fitness values. Otherwise, if the surrogates produce drastically different fitness values, the test case is re-evaluated on the SUT.

**Subject.** It refers to the artifacts approximated by a surrogate model. We identify two kinds of artifacts to which surrogates can be applied: *model* and *model&requirement assessment approximation*.

- *Model.* The model approximation technique relies on the creation of an approximation of the SUT that closely mimics its behavior, but is significantly cheaper to execute.

- *Model&requirement-assessment.* It relies on a function to predict whether test cases satisfy or violate the system requirements. That is, the surrogate approximates both the dynamics of the SUT and the assessment of the requirement.

**Design.** It refers to the process used to build a surrogate model. We identify two processes for the design of surrogate models: *manual* and *automatic*.

- *Manual.* In the manual process, engineers define a surrogate model by relying on their own experience and knowledge of SUT.

- *Automatic.* In the automatic process, engineers rely on some measurements to automatically construct a surrogate model. For example, these measurements can include logs recording the inputs and outputs of the SUT. We consider two variations for the automatic process: *dynamic* and *static*. The dynamic process improves surrogate models over time (e.g., thorough an approximation-refinement loop). The static process does not improve or modify surrogate models after they are computed.

**Type.** It captures the characteristics of surrogate models, and has two dimensions: the *Model type* and the *Output type.*

- *Model Type.* It captures the degree of information we have about the structure of surrogate models which can be one of the following three options: *white box*, *grey box*, and *black box.* The white box type indicates that surrogates are designed manually based on the knowledge that engineers have about the SUT. Hence, there is full observability to their structure. Black box models are used when the structure of the surrogate model is unknown, and engineers do not have any knowledge about the SUT, such as the physics and control laws that regulate the system behavior. Gray box models are used when there is partial information about the structure of surrogate models. For example, gray box models are useful when the engineer knows the physics and control laws that regulate the system behavior but does not know some parameters, e.g., the heat transfer coefficient of a car engine. Automatic techniques are needed to compute black box models and unknown parameters of grey box models.

- *Output Type.* It refers to the format of the outputs of surrogate models. We recognize two types of outputs: *Single Data* and *Time Series.* Some surrogate models produce a single data point as output, e.g., a Boolean value indicating whether a property is satisfied or a real value indicating a degree of satisfaction. Alternatively, surrogate models may generate time series data that is a sequence of data points indexed in time order.

**Context.** It refers to the context where surrogates are used. We identify three different usage contexts for surrogates: *Model-in-the-Loop (MiL)*, *Software-in-the-Loop (SiL)*, and *Hardware-in-the-Loop (HiL).*

- *Model-in-the-Loop.* It concerns with surrogate models that approximate a virtual and high-level representation of the system, e.g., some type of software model.

- *Software-in-the-Loop.* It relates to surrogate models that approximate software code.

- *Hardware-in-the-Loop.* It is about surrogate models that are built based on the measurements obtained from the actual operating system including its hardware components.

**Scope.** It refers to the portion of the search space surrogate models approximate. The search domain a surrogate model targets can be *global* or *local*.

- *Global.* When surrogate models capture a global view of the search space of the SUT, they approximate the entire search space, and hence, enable an explorative search.

- *Local.* It refers to the case in which surrogate models capture a local view by approximating a region of the search space that includes desired information. In that case, surrogates are utilized for an exploitative search.

**Evaluation Metrics.** It refers to the metrics used to assess an SBT approach. We identify three metrics: *Number of Iterations*, *Execution Time*, and *Estimated Execution Time*.

- *Number of Iterations.* It is the number of times an SBT approach executes the SUT.

- *Execution Time.* It measures the time the approach requires to detect a requirement violation.

- *Estimated Execution Time.* It predicts the time the approach would take to be applied to a compute-intensive system based on the results obtained for non-compute-intensive systems.

We will refect on these metrics in detail in Section 5.

## 3. Positioning against the State of the Art and Practice

To position the existing literature against the state of the arts and practice, one of the authors collected existing papers and classified relevant research literature along our taxonomy categories. We considered 14 papers, the four subject papers of our study Matinnejad et al. (2014); Abdessalem et al. (2016, 2018a); Menghi et al. (2020) and a set of ten papers by other researchers that employ SBT for testing complex software systems identified by forward snowballing Beglerovic et al. (2017); Arrieta et al. (2017); Wang et al. (2022); Innes and Ramamoorthy (2022); Zhang and Arcaini (2021); Pedrielli et al. (2021); Humeniuk et al. (2021, 2022); Zhong et al. (2021);

Haq et al. (2022), and classified relevant research literature along our taxonomy categories. We relied on Google Scholar for the forward snowballing: we identified new papers by consulting papers citing the four subject papers, we reviewed the retrieved papers, and we selected the papers that proposed a new SBT approach. The forward snowballing activity lead to ten papers, including one paper Zhong et al. (2021) proposing fuzz testing techniques and six papers presenting falsification techniques Innes and Ramamoorthy (2022); Zhang and Arcaini (2021); Pedrielli et al. (2021); Humeniuk et al. (2021, 2022). These works are specific instances of SBT techniques: fuzz testing is an SBT technique that injects unexpected inputs into a system to reveal software defects, and falsification techniques are SBT techniques driven by a requirement expressed using a formal language. We used these papers to demonstrate example research for different aspects of our taxonomy, and also, to analyze what sub-categories of our taxonomy are lesser studied compared to the others. Table 1 classifies relevant research literature along our taxonomy categories. Note that some papers are listed for some categories in Table 1. For example, for the evaluation metrics category, we do not list papers that do not evaluate their solutions. Table 1 shows that the SA-SBT techniques we reviewed cover all the nine top-level categories, and only two sub-categories (the Automatic-Static sub-category of the design category, and the SiL sub-category of the Context category) are not covered. In the following, we discuss the main findings for each category.

**Purpose**. Most of the techniques we reviewed use surrogates for fitness computation (8 out of 14), followed by test generation guidance (3 out of 14) and individual selection (3 out of 14). Although research literature has focused less on using surrogates for test generation guidance and individual selection compared to fitness computation, incorporating surrogates can greatly enhance SBT techniques by pinpointing critical regions of the test input domain and selecting the most crucial individuals (test cases).

**Usage**. Most of the techniques we reviewed use test case generation using single surrogates (7 out of 14), followed by test case generation using ensemble surrogates (4 out of 14), and test case filtering using single surrogates (3 out of 14). Although techniques that use single surrogates for test case generation are prevalent, surrogate models are also used for other testing activities (i.e., test case filtering) and aggregated into ensembles. Since surrogate models are a relatively new technology in SBT, their potential usages are still unclear and will grow over time.

Table 1: Classifying the SA-SBT techniques based on our taxonomy

| Category | References |
|---|---|
| **Purpose** | **Fitness Computation** Abdessalem et al. (2016); Matinnejad et al. (2014); Haq et al. (2022); Beglerovic et al. (2017); Arrieta et al. (2017); Pedrielli et al. (2021); Humeniuk et al. (2022); Wang et al. (2022) |
| | **Test GenerationGuidance** Abdessalem et al. (2018a); Humeniuk et al. (2021); Zhang and Arcaini (2021) |
| | **Individual Selection** Menghi et al. (2020); Zhong et al. (2021); Innes and Ramamoorthy (2022) |
| **Usage** | **TC generation–single** Menghi et al. (2020); Abdessalem et al. (2018a); Beglerovic et al. (2017); Arrieta et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Innes and Ramamoorthy (2022) |
| | **TC filtering** Abdessalem et al. (2016); Matinnejad et al. (2014); Zhang and Arcaini (2021) |
| | **TC generation–ensemble** Haq et al. (2022); Humeniuk et al. (2021, 2022); Wang et al. (2022) |
| **Subject** | **Model** Menghi et al. (2020); Beglerovic et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Zhang and Arcaini (2021); Innes and Ramamoorthy (2022) |
| | **Model and Requirement-Assessment** Haq et al. (2022); Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Arrieta et al. (2017); Humeniuk et al. (2021, 2022); Wang et al. (2022) |
| **Design** | **Manual** Arrieta et al. (2017); Humeniuk et al. (2021) |
| | **Automatic–Static** |
| | **Automatic–Dynamic** Menghi et al. (2020); Haq et al. (2022); Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Beglerovic et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Humeniuk et al. (2022); Zhang and Arcaini (2021); Wang et al. (2022); Innes and Ramamoorthy (2022) |
| **Model Type** | **White Box** Menghi et al. (2020); Beglerovic et al. (2017); Humeniuk et al. (2021, 2022); Zhang and Arcaini (2021); Innes and Ramamoorthy (2022) |
| | **Gray Box** Menghi et al. (2020); Humeniuk et al. (2022) |
| | **Black Box** Menghi et al. (2020); Haq et al. (2022); Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Arrieta et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Humeniuk et al. (2022); Wang et al. (2022) |
| **Output Type** | **Single Data** Haq et al. (2022); Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Arrieta et al. (2017); Zhong et al. (2021); Humeniuk et al. (2021, 2022); Wang et al. (2022) |
| | **Time Series** Menghi et al. (2020); Beglerovic et al. (2017); Pedrielli et al. (2021); Zhang and Arcaini (2021); Innes and Ramamoorthy (2022) |
| **Context** | **MiL** Menghi et al. (2020); Haq et al. (2022); Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Beglerovic et al. (2017); Arrieta et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Humeniuk et al. (2021, 2022); Zhang and Arcaini (2021); Wang et al. (2022); Innes and Ramamoorthy (2022) |
| | **SiL** |
| | **HiL** Beglerovic et al. (2017) |
| **Scope** | **Global Search** Menghi et al. (2020); Haq et al. (2022); Abdessalem et al. (2016, 2018a); Beglerovic et al. (2017); Arrieta et al. (2017); Zhong et al. (2021); Pedrielli et al. (2021); Humeniuk et al. (2021, 2022); Zhang and Arcaini (2021); Wang et al. (2022); Innes and Ramamoorthy (2022) |
| | **Local Search** Haq et al. (2022); Matinnejad et al. (2014); Pedrielli et al. (2021); Wang et al. (2022) |
| **Metric** | **Estimated Execution Time** Menghi et al. (2020) |
| | **Execution Time** Abdessalem et al. (2016); Matinnejad et al. (2014); Haq et al. (2022); Abdessalem et al. (2018a); Beglerovic et al. (2017); Arrieta et al. (2017); Humeniuk et al. (2021, 2022) |
| | **Num of Iterations** Menghi et al. (2020); Zhong et al. (2021) |

**Subject**. For 6 out of 14 papers, the subject of the surrogate model is the original model under test, while for the remaining eight papers, it includes both the model and the requirement assessment. Using a surrogate model to approximate the original model under test and a fitness function to evaluate the different test cases ensures that the system's dynamics are approximated while the fitness computation is accurate. Differently, also considering an approximation of the fitness function introduces additional noise. There is a need for more empirical studies to precisely identify the situations where one of these approaches is preferred.

**Design**. Most of the techniques we reviewed (12 out of 14) automatically compute the surrogate model and use dynamic processes that improve the surrogate model over time. Two techniques manually defined the surrogate model. None of the techniques we reviewed used a static surrogate model, i.e., a fixed surrogate model that is not improved over time. Although automatic techniques do not require manually defining the surrogate model, in many applications, the manual design offers several benefits since engineers can incorporate knowledge about the structure of the SUT in the design of the surrogates. More empirical studies are needed to identify the scenarios in which one technique must be preferred over the other.

**Model Type**. Most of the techniques we reviewed are black box (8 out of 14), followed by white box (4 out of 14). Two techniques, Menghi et al. (2020); Humeniuk et al. (2022), are generic and rely on other technologies such as system identification (SI) Bittanti (2019); AMS and Electromagnetism (1989) that enable them to support black box, white box, and gray box model types. In particular, SI enables us to create surrogates that belong to all the three categories. Although less studied, gray box technologies can offer a good trade-off in situations in which engineers have some knowledge about the system's internal structure and need to be aware of the values of some of the parameters that still need to be discovered.

**Output Type**. For 9 out of 14 techniques, the output type is Single Data. The remaining five techniques use time series. Therefore, in most cases, the surrogates directly predict the single data fitness value for a test input. Surrogates that use time series as output type are more recent Menghi et al. (2020). Time series data leads to the generation of more data points that help with better training of surrogate models.

**Context**. For all the techniques we considered, the surrogate context is MiL. One technique Beglerovic et al. (2017) claims that their approach also supports HiL in addition to MiL. The prevalence of MiL techniques is justified by the fact that MiL relies on simulators and system models, which are often publicly available. In contrast, SiL and HiL involve running the software under test on real or emulated hardware and require direct access to the source code and system hardware, which can be challenging in an industrial context. To overcome this hurdle, close collaboration with industry partners is often necessary. Furthermore, creating the necessary hardware infrastructure and setup can be both costly and time-consuming, particularly in an academic setting.

**Scope**. The scope of 10 out of 14 techniques is global search. The scope of one techniques is exclusively local search. Three techniques consider both for local and global search. We noticed that the boundary between global and local search was often blurry. Additional research work is needed to precisely identify and define the boundary between the local and global search for SBT.

**Evaluation Metric**. Most of the techniques use execution time as an evaluation metric (8 out of 14), followed by the number of iterations (2 out of 14), and the estimated execution time (1 out of 14). Some techniques either did not use any evaluation metric, as they propose theoretical solutions which are not empirically evaluated, or their evaluation did not involve a comparison of SA-SBT solutions. Although the execution time measures the actual time required for fault-finding, it makes experiments hardly reproducible as it depends from the actual hardware platform on which the experiments are executed. On the contrary, the number of iterations enables experiment replication as the hardware platform does not influence it. However, the execution time of the iterations may differ across different techniques. We provide additional reflections and considerations on the evaluation metric category in Section 5.

## 4. Replication Studies

Papers listed in Table 1 do not present studies that try to replicate the approaches of the subject papers. In the next two sections, we present two replication studies based on two of the subject papers. Specifically, we report replication studies of NSGAII-DT Abdessalem et al. (2018a) and ARIsTEO Menghi et al. (2020). NSGAII-DT and ARIsTEO belong to different sub-categories under the purpose, subject, model type, output type, and evaluation metric categories in our taxonomy. Replicating these complementary approaches allows us to gain insightful observations and lessons learned.

### 4.1. Replication Study: ADAS Case Study

In this section, we report on the replication of the experiments from the subject paper proposing the NSGAII-DT technique, which was evaluated using an industrial ADAS case study Abdessalem et al. (2018a). The replication attempt was performed by different researchers (the second and third authors of the paper) and by considering an ADAS which has been developed in collaboration with industrial partners and is different from that used in
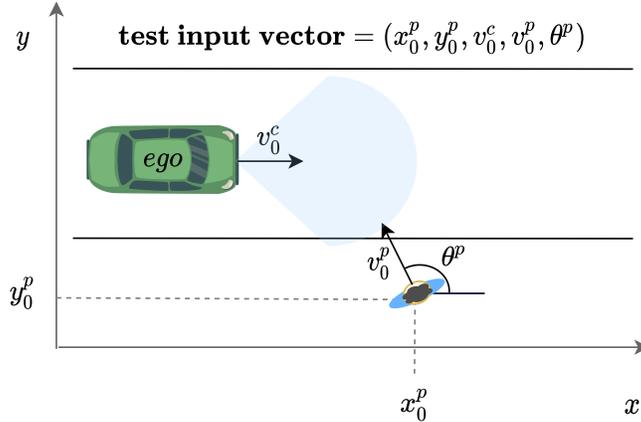
Figure 3: SUT of the original study Abdessalem et al. (2018a): a pedestrian crosses the lane of the ego car equipped with the AEB system.

the original paper Abdessalem et al. (2018a). In the following, we first describe the experiment setup of the original and the replication studies, then we present our results.

### 4.1.1. Experiment setup

**System under test (SUT)**. The system under test for the original subject paper is the Automated Emergency Breaking (AEB) system presented in Figure 3. The AEB identifies pedestrians in front of a vehicle and avoids collision by applying the brake when necessary. The ego car (i.e., the self-driving car equipped with AEB) drives from the left-end on an urban street, and a pedestrian starts crossing the street from the middle of the street.

For the replication study, we consider an Automated Valet Parking (AVP) system presented in Figure 4: A driver leaves the car at the entrance of a parking lot and AVP autonomously parks the car in an available parking spot. The ego car is equipped with a Lidar sensor and an AEB to ensure that the car does not collide with other objects, e.g., pedestrians or cars. The AVP is tested using a scenario where an occluded pedestrian appears in the path of the ego car as the ego car approaches a free parking spot. An occluded pedestrian is a pedestrian who is partially or completely obscured from the ego car's view by another object (e.g., a pillar or a parked vehicle).

**Search Space**. The original study tests are obtained by varying the following features of the AEB test scenario (see Figure 3): The speed of the

15

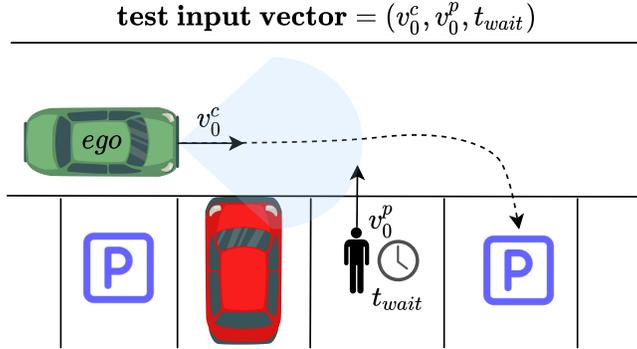**test input vector** $= (v_0^c, v_0^p, t_{wait})$

Figure 4: AVP testing scenario from the replication study: a pedestrian occluded by a parking car is crossing the lane of the ego car, which is approaching a free parking spot.

ego car $(v_0^c)$, the initial position of the pedestrian with respect to the ego car $(x_0^p, y_0^p)$, the speed $(v_0^p)$ and the orientation $(\theta^p)$ of the pedestrian, the weather type, which can be either normal, rainy and snowy, and the road shape, which can be either straight, curved or ramped. For the curved and ramped roads, the test input includes the curve radius and the ramp height. For the snowy and rainy weather types, the test input includes the level of precipitation and the presence of fog with different density levels.

The search space for the replication study varies the speed of the ego car $(v_0^c)$, the speed of the pedestrian $(v_0^p)$ and the time that the pedestrian starts running after the start of the simulation $(t_{wait})$. The variable $t_{wait}$ has a similar purpose as the variables $(x_0^p, y_0^p)$ in the original study. Both $t_{wait}$ and $(x_0^p, y_0^p)$ help generate different tests where the ego car and the pedestrian appear in different proximities to one another. The AVP testing excludes road shape and weather-related factors, as they are not essential for assessing AVP functionality. Given the confined parking lot environment, fluctuations in these conditions are considered non-essential for evaluating system performance. In addition, the AVP testing scenario is focused on occluded pedestrians, and for them, the initial position and the orientation are not relevant since they are blocked from the ego car's field of view at the beginning. Instead, for occluded pedestrians, the relevant information for testing is when they appear in front of the ego car (variable $t_{wait}$), which is included as a test input in the replication study.

Both the original and the replication studies use the PreScan simulator pre (2023) to test their respective SUT. Both studies set 100 as the sampling rate for PreScan and generate simulations with the duration of 10s.

**Fitness functions.** Both studies define multiple fitness functions to assess the quality of individual simulations. Fitness functions measure how close the scenario is to a collision. The original study defines three fitness functions: ($F_1$) The minimum distance between the pedestrian and the field of view of the ego car, ($F_2$) The velocity of the ego vehicle at the time where the distance between the ego car and the pedestrian is at its lowest value, i.e., when the minimum for $F_1$ is computed, and ($F_3$) The certainty of the detection of a pedestrian. By minimizing $F_1$ and maximizing $F_2$, we obtain scenarios that either represent near-collision situations or collisions with a high speed of the ego car. In addition, the original study considers maximising $F_3$ since the study was focused on identifying scenarios where AEB detects a pedestrians but fails to avoid a collision with the pedestrian. Otherwise, scenarios where the pedestrian cannot be detected are due to the failure of the object detection component and not the failure of AEB.

The replication study defines two fitness functions: (1) $F_1'$ (similar to $F_1$) this function computes the minimum distance between the ego car and the pedestrian. However, in the replication study, this distance is defined to be the distance between the pedestrian and the front side of the car. Hence, the $F_1'$ function is defined such that it assigns more optimal values to collisions where the pedestrian hits the car from the front. This is to provide more effective guidance for the search to generate scenarios where the collision is realistic and meaningful. That is, the front side of the ego car hits a pedestrian. (2) $F_2'$ (similar to $F_2$) this function captures the velocity of the ego vehicle at the time where the distance between the ego car and the pedestrian is at its lowest value. The replication study does not consider the third fitness function (i.e., the certainty of detection) as in the original study, because the AVP does not provide detailed output that measures the certainty of detection.

In addition to fitness functions, both studies define when a scenario represents a failure of the SUT. For both studies, a failure is when the ego car hits the pedestrian with a high speed. Failure revealing scenarios are identified by setting a threshold on $F_1$ and $F_2$ (resp. $F_1'$ and $F_2'$ in the replication study) to indicate that the car and the pedestrian collide and the car has a high speed at the time of the collision. Since the SUTs in the original and the replication studies are different, the threshold values for what defines a

failure in these two studies are different too. Hence, the thresholds in these studies were defined independently and based on the requirements of their respective SUTs.

**Initial population**. The replication study uses Latin Hypercube Sampling McKay et al. (1979), a quasi-random sampling strategy, to generate the initial population. The original study used combinatorial testing using the Pledge tool Henard et al. (2013) to generate an initial population of diversified individuals. Given that the purpose of both strategies is to generate a diverse initial set of individuals, we believe the differences in generating the initial populations do not significantly impact the results of our study.

**Search Algorithm**. Both the original and the replication study rely on NSGAII-DT. NSGAII-DT uses machine learning decision tree models and multi-objective evolutionary search Luke (2013). Decision tree models are developed based on the scenarios generated at intermediary search iterations. The models learn the characteristics of the critical test scenarios and identify critical regions in an input space (i.e., the regions of a test input space that are likely to contain more fault-revealing test cases and test cases that reveal more severe faults). The subsequent search iterations then focus on the critical regions, and select and evolve critical scenarios in those regions. When the search inside the critical regions is terminated, a new decision tree model is re-generated using the entire population. NSGAII-DT continues by iteratively building decision trees followed by search iterations focused on critical regions identified by the trees until the search time budget runs out. The original and replication studies consider the same values for the search parameters, such as mutation and crossover rates. The main parameters and setup elements of the replication and the original studies are listed in Table 2.

### 4.1.2. Results

The original study provides two sets of metrics to evaluate the results of NSGAII-DT: (1) Quality indicators for Pareto-based search algorithms Chen et al. (2020), and (2) the number of distinct critical scenarios representing failures of the SUT. As a Pareto-based search algorithm, NSGAII-DT aims to find a set of non-dominated solutions in a multi-objective optimization problem. These algorithms are typically assessed based on a special category of metrics, referred to as *quality indicators* Chen et al. (2020) that compare Pareto-based search algorithm in terms of convergence, uniformity, spread and cardinality. In particular, the original study used three well-known quality indicators, Hypervolume (HV), Spread ($\Delta$) and Generational

Table 2: Experimental setup for replication study vs. the study in original paper  Abdessalem et al. (2018a)

| Category | Replication study setup | Original study setup |
|---|---|---|
| SUT | Longitudinal/Latitudinal vehicle dynamics, Path Follower, Lidar | 3D vehicle dynamics, Object detection, Lidar |
| Sampling rate | 1/100 | 1/100 |
| Scenario | Pedestrian crossing street | Pedestrian crossing street |
| # fitness functions | 2 | 3 |
| Initial sampling | LHS | Combinatorial Testing |
| Mutation probability | 1/n | 1/n |
| Crossover probability | 0.6 | 0.6 |

Distance (GD), to compare the results of NSGAII-DT with those of NSGAII. Below we describe the experiment setup for our replication study and then compare and discuss the results of both studies.

**Setup.** We evaluate the performance of NSGAII-DT and NSGAII on the AVP case study for a minimum of 300 minutes. We note that NSGAII served as a baseline in the original study. The average execution time of a single AVP simulation is approximately 19 seconds, allowing for at least 1000 simulations (fitness evaluations) to be completed within the 300-minute computation time. We rerun each of NSGAII-DT and NSGAII for 10 times to account for their randomness. The experiments were performed on a computer with an i7-8700 processor with 3.19 Ghz (12 cores) and 32 GB RAM.

Note that the AEB simulations in the original study took significantly longer than the AVP simulations, with an average execution time of 1.2 minutes for AEB and 19 seconds for AVP. In the original study, both NSGAII-DT and NSGAII algorithms were run for a 24-hour period, providing enough time for at least 1000 minimum AEB simulations, similar to our replication study.

**Metrics.** For NSGAII, we compute HV, Spread ($\Delta$) and GD after each generation. But, for NSGAII-DT, we compute these metrics after each generation of each of the NSGAII algorithm runs performed inside the critical regions. HV measures the size of the space covered by the members of a Pareto front generated by a search algorithm Chen et al. (2020). The higher this size, the better the results of the algorithm. GD measures the Euclidean

distance between members of the calculated Pareto front and the nearest solutions on a reference Pareto front Chen et al. (2020). The lower the value of GD, the more optimal the calculated Pareto front solutions are. Spread measures the extent of spread among the members of a Pareto front generated by a search algorithm Deb et al. (2002). The lower the Spread values, the better spread out the search outputs. In addition, we compute the number of distinct critical scenarios generated by NSGAII and NSGAII-DT, following the definition provided in the original study, where distinct scenarios are defined as those that differ in at least one input value Abdessalem et al. (2018a).

**Data Availability.** The implementation of NSGAII and NSGAII-DT for AVP and the detailed results of these algorithms for our replication study are available online res (2023).

**Comparison of the results.** Figures 5 and 6 show the results of the HV, Spread and GD metrics over time for NSGAII and NSGAII-DT for our original and replication study, respectively. The results for the original study are shown at every four-hour time interval starting at 2h until the time limit of 24h, while the results for the replication study are shown at every 60 min starting from 60min and until 300min. Since the original and replication experiments are performed on different case studies with different inputs and fitness functions, we cannot make direct comparisons between the values of these metrics, but instead, we aim to compare the trends of these metrics over time.

In both the original and replication study, the HV values for NSGAII-DT are better than those of NSGAII at the start of the experiments. However, over time, the differences between the HV values of the two algorithms tend to decrease. In the replication study, this reduction is quite steep, with the HV values of NSGAII and NSGAII-DT eventually becoming equal. In contrast, in the original study, the HV values of the two algorithms seem to stabilize at a certain point. The GD values of NSGAII-DT are slightly better than those of NSGAII over the entire duration both in the original study and in the replication study. However, at the end of the simulation, which is 300 minutes, the GD values of NSGAII slightly improve over those of NSGAII-DT in the replication study. Further, the variations of the GD values for NSGAII-DT in the replication study are noticeably higher than those of the GD values of NSGAII-DT in the original study. There are remarkable differences between the replication and the original studies in the values of the Spread metric obtained for NSGAII-DT. While in the original study the values of Spread for NSGAII-DT are considerably better than those
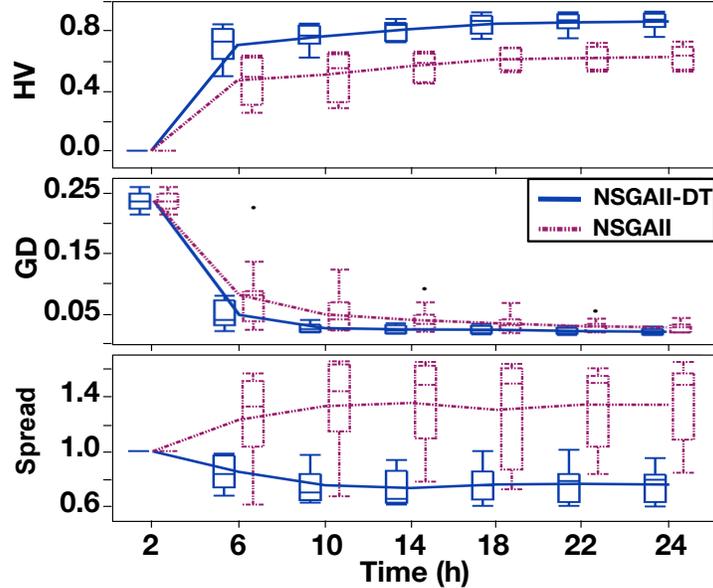
Figure 5: The comparison of HV, GD and Spread (Δ) values obtained by NSGAII and NSGAII-DT from the original study Abdessalem et al. (2018a).

of NSGAII, for the replication study the values of Spread for both NSGAII and NSGAII-DT are almost the same.

Table 3 shows the number of distinct critical scenarios generated by NSGAII-DT and NSGAII in the replication and original studies. As shown in the table, in both studies, provided with the same search time budget, NSGAII-DT is able to find considerably more distinct critical scenarios compared to NSGAII. In particular, the number of distinct critical scenarios that NSGAII-DT finds are, respectively, 1.8 and 2.7 times more than those that NSGAII can find in the original and replication studies.

**Summary.** The replication study supports some of the conclusions of the original study. However, there are some disparities with some specific findings. Both studies confirm that the surrogate-based SBT method, NSGAII-DT, generates significantly more distinct failures than its non-surrogate counterpart within the same search time budget. However, the consistency of the results comparing the quality of the Pareto-based solutions generated by NSGAII and NSGAII-DT is questionable. The HV and GD results from the replication study partially align with those from the original study, whereas
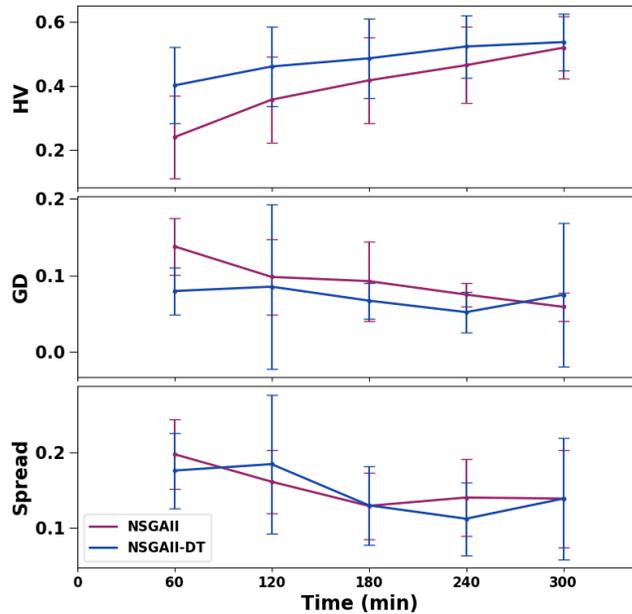
Figure 6: The comparison of HV, GD and Spread values obtained by NSGAII and NSGAII-DT from the replication study.

the results of the Spread metric are significantly different. Despite our efforts to align the two studies as much as possible, the SUTs, and hence, the test inputs and fitness functions remain to be different. These factors may have contributed to the discrepancies in the outcomes. Independently from the real sources of these differences, our replication study allowed us to gain a deeper understanding of the limitations and significance of the metrics used to evaluate NSGAII-DT, which we elaborate on in Section 5.

## 4.2. Replication Study: Simulink

In this section, we report on an attempt to replicate the experiments of ARIsTEO, an approximation-refinement testing technique driven by surrogate models proposed by one of the subject papers Menghi et al. (2020). The ARIsTEO technique is showcased by a plugin of S-TaLiRo Annpureddy et al. (2011), a falsification-based testing tool for Simulink models. In their

Table 3: The number of distinct critical scenarios obtained by NSGAII and NSGAII-DT from the replication and original study.

|  | Replication study | Original study |
|---|---|---|
| NSGAII | 639 | 411 |
| NSGAII-DT | 1733 | 731 |
| Ratio NSGAII-DT / NSGAII | 2.7 | 1.8 |

paper, the authors of ARIsTEO, evaluated their contribution by considering a set of case studies which included some of the benchmark models of the ARCH competition G. Ernst et al. (2021) and one additional industrial case study from the satellite domain. The industrial case study is not publicly available since the authors could not release it due to a non-disclosure agreement. However, the authors provided a complete replication package with the other benchmark models and results. In addition, they participated in the ARCH competition, editions of 2020 Ernst et al. (2020), 2021 G. Ernst et al. (2021), and 2022 Ernst et al. (2022). The ARCH competition is a friendly yearly competition between testing tools for continuous and hybrid systems ARCH. The competition includes several participants, such as Fal-CAuN Waga (2020), S-TaLiRo Annpureddy et al. (2011), falsify Yamagata et al. (2021), FalStar Ernst et al. (2018), ForeSee Zhang et al. (2021). The models and requirements considered by the competition are publicly available and provided under different licenses.

We report on our attempt to replicate the results generated from ARIsTEO from editions 2020 Ernst et al. (2020), 2021 G. Ernst et al. (2021), and 2022 Ernst et al. (2022).

*4.2.1. Experiment Setup*

Our benchmark consists of seven models: Automatic Transmission (AT) Hoxha et al. (2015), Fuel Control on Automotive Powertrain (AFC) Jin et al. (2014), Neural Network Controller (NN)NN (2022), Wind Turbine (WT) Schuler et al. (2017), Chasing Cars (CC) Hu et al. (2000), Aircraft Ground Collision Avoidance System (F16) Heidlauf et al. (2018), and Steam Condenser (SC) Yaghoubi and Fainekos (2019). These Simulink models have different sizes and complexity: the number of Simulink blocks spans from 13 (CC) to 302 (AFC). The models are from different domains: automotive (AT, AFC), neural networks (NN), and energy (SC). AFC is from Toyota. The models come with different licenses, including GPLv2.

Table 4: Results for piecewise continuous input signals.

| | 2020 | | | 2021 | | | 2022 | | |
|---|---|---|---|---|---|---|---|---|---|
| | $FR$ | $\overline{S}$ | $\widetilde{S}$ | $FR$ | $\overline{S}$ | $\widetilde{S}$ | $FR$ | $\overline{S}$ | $\widetilde{S}$ |
| AT1 | 0 | - | - | 0 | - | - | 0 | - | - |
| AT2 | 50 | 4.4 | 4 | 50 | 4.1 | 3 | 10 | 5.4 | 2.0 |
| AT51 | | nSbyA | | 0 | - | - | 0 | - | - |
| AT52 | | nSbyA | | 50 | 3.2 | 2 | 10 | 1.7 | 1.5 |
| AT53 | | nSbyA | | 50 | 2.6 | 2 | 10 | 12.8 | 8.5 |
| AT54 | | nSbyA | | 3 | 295.1 | 300 | 0 | - | - |
| AT6a | 45 | 90.7 | 69 | 50 | 39.0 | 23 | 10 | 93.4 | 70.0 |
| AT6b | 50 | 18.1 | 15 | 49 | 101.7 | 83 | 4 | 112.3 | 127 |
| AT6c | 44 | 95.6 | 66 | 50 | 27.2 | 16 | 10 | 124.2 | 111.5 |
| AT6abc | | NPofB | | 50 | 27.2 | 16 | 10 | 73.4 | 30.5 |
| NN | 50 | 62.8 | 46 | 50 | 62.8 | 46 | 1 | 299.0 | 299.0 |
| NNx | | NPofB | | | NPofB | | 0 | - | - |
| WT1 | 50 | 15.6 | 10 | 50 | 13.0 | 10 | | NPofB | |
| WT2 | 50 | 1.5 | 1 | 50 | 1.3 | 1 | | NPofB | |
| WT3 | 50 | 3.4 | 3 | 50 | 2.4 | 2 | | NPofB | |
| WT4 | 50 | 1.0 | 1 | 50 | 1.0 | 1 | | NPofB | |
| CC1 | 50 | 16.1 | 11 | 50 | 27.0 | 15 | 10 | 24.8 | 17.5 |
| CC2 | 50 | 1.0 | 1 | 50 | 9.2 | 6 | 10 | 14.6 | 9.0 |
| CC3 | 50 | 45.8 | 27 | 50 | 56.3 | 46 | 10 | 29.5 | 35.0 |
| CC4 | 50 | 1.0 | 1 | 0 | - | - | 0 | - | - |
| CC5 | 49 | 52.5 | 40 | 50 | 25.9 | 17 | 10 | 18.8 | 18.0 |
| CCx | | NPofB | | 15 | 250.0 | 300 | 4 | 134.0 | 74.0 |
| F16 | | nSbyA | | 0 | - | - | 0 | - | - |
| SC | 50 | 1.0 | 1 | 0 | - | - | 0 | - | - |

The models of the ARCH competition are associated with 27 requirements expressed using a logic-based language. Each requirement is associated with an identifier that starts with the model's name. The rows of Table 4 and Table 5 represent different requirements. For example, the row with the identifier AT6a refers to the requirement AT6a of the AT model. The interested reader can find additional information about the models and requirements in the report of the ARCH competition (i.e., G. Ernst et al. (2021); Ernst et al. (2020, 2022)).

The same configuration is considered for ARIsTEO across the three years it participated in the competition. Specifically, ARIsTEO is configured as fol-

Table 5: Results for constrained input signals.

| | 2020 | | | 2021 | | | 2022 | | |
|---|---|---|---|---|---|---|---|---|---|
| | $FR$ | $\overline{S}$ | $\widetilde{S}$ | $FR$ | $\overline{S}$ | $\widetilde{S}$ | $FR$ | $\overline{S}$ | $\widetilde{S}$ |
| AT1 | 0 | - | - | 0 | - | - | 0 | - | - |
| AT2 | 50 | 4.5 | 3 | 50 | 5.1 | 4 | 10 | 12.9 | 8.0 |
| AT51 | | nSbyA | | 50 | 8.4 | 6 | 10 | 19.0 | 10.0 |
| AT52 | | nSbyA | | 50 | 3.5 | 3 | 9 | 74.7 | 46.0 |
| AT53 | | nSbyA | | 50 | 2.6 | 2 | 10 | 1.4 | 1.0 |
| AT54 | | nSbyA | | 50 | 29.4 | 17 | 10 | 44.0 | 40.0 |
| AT6a | 41 | 116.3 | 72 | 47 | 103.1 | 80 | 7 | 65.1 | 62.0 |
| AT6b | 50 | 36.3 | 27 | 45 | 164.7 | 174 | 5 | 91.2 | 74.0 |
| AT6c | 44 | 89.8.6 | 73 | 49 | 89.1 | 60 | 7 | 175 | 141.0 |
| AT6abc | | NPofB | | 50 | 79.1 | 70 | 9 | 84.1 | 81.0 |
| AFC27 | 50 | 2.3 | 1 | 50 | 2.3 | 1 | 9 | 21.9 | 25.0 |
| AFC29 | 50 | 28.5 | 23 | 50 | 1.0 | 1 | 10 | 3.4 | 3.5 |
| AFC33 | 50 | 24.7 | 16 | 50 | 1.0 | 1 | 0 | - | - |
| NN | 50 | 62.8 | 46 | 50 | 62.8 | 46 | 10 | 117 | 82.5 |
| NNx | | NPofB | | 50 | 1.0 | 1 | 0 | - | - |
| WT1 | 50 | 1.4 | 1 | 50 | 1.4 | 1 | | NPofB | |
| WT2 | 50 | 1.0 | 1 | 50 | 1.0 | 1 | | NPofB | |
| WT3 | 50 | 1.1 | 1 | 50 | 1.1 | 1 | | NPofB | |
| WT4 | 50 | 1.0 | 1 | 50 | 1.0 | 1 | | NPofB | |
| CC1 | 50 | 28.8 | 22 | 50 | 9.1 | 6 | 10 | 9.1 | 8.0 |
| CC2 | 50 | 1.0 | 1 | 50 | 6.7 | 4 | 10 | 10.8 | 9.0 |
| CC3 | 50 | 18.1 | 16 | 50 | 18.9 | 15 | 10 | 12.8 | 13.4 |
| CC4 | 50 | 1.0 | 1 | 0 | - | - | 0 | - | - |
| CC5 | 48 | 66.9 | 44 | 50 | 29.1 | 12 | 10 | 21.1 | 11.5 |
| CCx | | NPofB | | 20 | 229.4 | 12 | 3 | 97 | 103.0 |
| SC | 50 | 1 | 1 | 0 | - | - | 0 | - | - |

lows. We used an ARX model (ARX-2) with order $na = 2$, $nb = 2$, and $nk = 2$[1] as the structure for the surrogate model used in the approximation-refinement loop of ARIsTEO. For models with multiple inputs and outputs, the dimension of the matrix $na$, $nb$, and $nk$ are changed depending on the number of inputs and outputs. For the other parameters, we assigned the default configuration values of S-TaLiRo. We considered the same parametrization of S-TaLiRo for the input signals. ARIsTEO executes the original Simulink

---

[1] https://nl.mathworks.com/help/ident/ref/arx.html

model to learn the initial surrogate model. The competition sets the cut-off values for the number of simulations of the original model to 300 for 2020 and 2021 and do not impose a maximal number of simulations that can be run in one falsification trial for 2022. This choice was performed to enable a more accurate assessment of the difficulty of the benchmarks. We considered this value also for the number of simulations of the surrogate model (per trial). However, we considered 300 as the maximum number of iterations for ARIsTEO in 2022 since we aimed to maintain consistency with the previous years of the competition. Notice that this choice penalizes ARIsTEO over the other competing tools in 2022.

During the competition, the participants must test the model for each requirement: *piece-wise continuous input signals*, and *constrained input signals*. The first option let the participant decide the shape of the input signals, and the second option precisely fixes the format of the input signal. Table 4 and Table 5 respectively report the results obtained by ARIsTEO for the two options. In addition, the competition required participants to execute experiments 50 times in 2020 and 2021 and 10 times in 2022 due to the non-determinism of the underlying search-based algorithms. In 2022 the number of experiments was reduced from 50 to 10 since there was no upper bound on the maximal number of simulations that can be run in one falsification trial.

*4.2.2. Results*

Table 4 and Table 5 report the results of our experiments. For each year, they report the falsification rate with respect to the number of trials ($FR$), mean ($\overline{S}$) and median ($\widetilde{S}$) number of simulations (rounded down) over successful trials ("–" if $FR$ is zero). Cells labeled with NPofB (Not Part of the Benchmark) specifies that the model was not part of the benchmark models considered by the competition that year. For the first year (2020) of the competition, the cells for the benchmarks AT51, AT52, AT53, and AT54 are labeled with nSbyA (not supported by ARIsTEO) since ARIsTEO could not process the benchmark. Indeed, for the first year of the competition, ARIsTEO could not support property specifications that use locations. Locations are used by S-TaLiRo Annpureddy et al. (2011) (the tool ARIsTEO extends) to specify properties that involve variables that can assume values within a finite set of values.

The results from Table 4 show that the falsification rates for AT1, AT2, AT51, AT52, AT53, AT6abc, WT1, WT2, WT3, WT4, CC1, CC2, CC3,

CC5, F16 are the same for all the three years of the competition. The mean ($\overline{S}$) and median ($\widetilde{S}$) number of simulations (rounded down) over successful trials are also consistent for these benchmarks (small fluctuations are caused by the non-determinism of the underlying search-based algorithm). For AT54, the property was not falsified in 2022 while it was falsified in 3 cases in 2021. However, unlike 2021 that require the experiments to be executed 50 times, in 2022 experiments were only executed 10 times. For AT6a, AT6b, AT6c, and CCX small fluctuations in the falsification rate, mean ($\overline{S}$) and median ($\widetilde{S}$) number of simulations are caused by the non-determinism of the underlying search-based algorithm which may produce different results across multiple experiments. For NN, we noticed that the results of 2022 are significantly different from the ones obtained in 2020 and 2021. Unlike 2022, in 2020 and 2021, there was a mistake in the specification of the requirement: there was a typo in one of the temporal operators used to specify the requirement ($\square_{[1,18]}$ instead of $\square_{[1,37]}$). For CC4 and SC, the results of 2020 (the first year in which ARIsTEO participated in the competition) are significantly different from the ones obtained in 2021 and 2022. For CC4, in 2020, the predicate was written incorrectly: $y_4 - y5 \geq 8$ instead of $y_5 - y4 \geq 8$. The former is immediately falsified. Furthermore, the simulation time in 2020 was 30s, which was lower than the temporal operators of the requirement (in 2021-22 it was 100s). For SC, in 2020, ARIsTEO used 12 control points and a simulation time of $100s$. In 2021-22, ARIsTEO used 20 control points, and the simulation time was reduced to $35s$. Furthermore, for piecewise continuous input signals ARIsTEO reports the results for constrained input signals: the piecewise continuous input signals let the participant be free to decide the shape of their input signals, and it was decided to use the one mandated by the constrained input signals.

The results from Table 5 show that the falsification rates for AT1, AT2, AT51, AT53, AT54, AFC29, NN, WT1, WT2, WT3, WT4, CC1, CC2, CC3, are the same for all the three years of the competition. The mean ($\overline{S}$) and median ($\widetilde{S}$) number of simulations (rounded down) over successful trials are also consistent for these benchmarks (small fluctuations are caused by the non-determinism of the underlying search-based algorithm). For AT52, AT6a, AT6b, AT6c, AT6abc, AFC27, and CC5 small fluctuations in the falsification rate, mean ($\overline{S}$) and median ($\widetilde{S}$) number of simulations are caused by the non-determinism of the underlying search-base algorithm which may produce different results across multiple runs. For AFC33, we noticed that the re-

sults of 2022 are significantly different from the ones obtained in 2020 and 2021. AFC33 has a different input range compared to AFC27 and AFC29. In 2020 and 2021, the input range of AFC27 and AFC29 was erroneously considered. In 2020, the wrong input range was considered (there was a bug in the code that was supposed to change the range). In 2021, the input range was correct, but the control points were considered in the wrong order. For NNx, the results of 2022 are significantly different from the ones obtained in 2021 (NNx was not part of the benchmark models in 2020). NNx has a different input range compared to NN. However, in 2021, the input range for NN was erroneously considered. For CC4, the results of 2020 are significantly different from the ones obtained in 2021 and 2022. As previously reported, the cause was a problem in the specification of one of the predicates and a different simulation time. For SC, the success rate for 2022 should be 0; there was a typo on the report for ARCH 2022. In addition, for SC, in 2020, ARIsTEO considered a different simulation time and a different number of control points.

## 5. Lessons Learned

In this section, we reflect on the lessons learned from our replication studies and the creation of a taxonomy for the SA-SBT literature. Our lessons concern the methods and metrics used to evaluate SA-SBT solutions, the scope and purpose of using surrogates in SBT, and the need for frameworks to facilitate extensive experimentation in this domain. We believe our lessons would be most relevant for software engineering researchers and practitioners working on the verification and testing of cyber-physical and autonomous systems.

### 5.1. Metrics for Evaluating SA-SBT

The evaluation metrics used in the subject papers and the research reviewed in Section 3 can be grouped into three categories: (1) *Metrics to assess the effectiveness of SA-SBT.* The effectiveness of a SA-SBT solution is evaluated by its ability to identify failures in the SUT and is quantified by the number of distinct failure scenarios it detects and the severity of each scenario, as indicated by the fitness function values. (2) *Metrics to assess the efficiency of SA-SBT.* As discussed in Section 2 under the *evaluation metrics* category, the efficiency of SA-SBT techniques is typically evaluated based on the number of iterations, execution time and estimated execution

time. (3) *Metrics to assess the performance SA-SBT solutions that rely on Pareto-based search algorithms.* These include metrics such as HV, Spread and GD discussed in Section 4.1. Below, we discuss the limitations of some of these metrics when applied to SA-SBT solutions, and suggest alternatives to these metrics or ways to redefine these metrics to be applicable to SA-SBT solutions.

**Metrics for assessing effectiveness.** Papers on software testing often report the number of generated failures as a metric to evaluate effectiveness (e.g., Abdessalem et al. (2018a); Haq et al. (2022)). However, for SA-SBT solutions, particularly when they are applied to case studies such as ADAS, it is important to report failures that are *distinct* and *valid*. Failures are considered distinct when they are generated by test inputs that are significantly different. In the study that we replicated in Section 4.1, distinct scenarios are defined as those that differ in at least one input value Abdessalem et al. (2018a). While this definition might be sufficient for the AEB and AVP case studies in Section 4.1, it may fall short in other situations where two test inputs may differ in variables that are not related to failures. For example, we may have two pedestrians in the initial scene where one is far away from the collision. Changing the position of the far-away pedestrian would not help reveal interesting and diverse failure situations. In addition, small variations in continuous input variables may not reveal diverse failures. For example, in Figure 3, slightly moving the pedestrian to the left or right in a failing test input may still lead to failures that are essentially the same. In a recent study, Zhong et. al. Zhong et al. (2021) proposed an improved definition of distinctness that requires variations in a minimum number of input variables and differences between continuous variables to exceed a user-defined threshold. However, this definition relies on user-defined parameters that can vary between domains, making it difficult to replicate and generalize the evaluation results.

In addition to being distinct, failures should be valid and meaningful. For example, both AVP and AEB fail if the car speed is higher than the speed limit for urban areas (i.e., $100km/h$). To ensure the validity of failures, it is important that we select the input variable ranges such that they satisfy the expected preconditions of the SUT. Unmet preconditions may cause failures that are not faults in the SUT, but rather caused by invalid inputs. The validity of failures, further, depends on the way we formalize the notion of failure or the way we define test oracles. In our NSGAII-DT study, the criticality function acts as a test oracle and determines whether, or not, a

29

scenario is a failure. For the AEB case study, the criticality function was defined using a conjunction of predicates over fitness functions that are capped by some thresholds. The formalization of the fitness functions, however, had limitations which led to some invalid failures, e.g., the pedestrian hitting the car from the side Abdessalem et al. (2016, 2018a). We addressed this issue, in our replication study (Section 4.1), by redefining the fitness functions such that scenarios where the pedestrians approach the ego car from the sides are penalized and are not considered as critical.

**Metrics for assessing efficiency.** A primary research question for evaluating SA-SBT techniques is to assess whether they *improve the efficiency of testing while ensuring its effectiveness.* All the SA-SBT techniques we reviewed in Table 1 answer this research question by comparing SA-SBT with some non-surrogate SBT baseline. They use different metrics for this comparison though. Specifically, all the approaches, except for Menghi et. al. Menghi et al. (2020), use the *execution time* metric. That is, they execute the SA-SBT and the baseline techniques for the same amount of time and compare the results. In all these cases, however, the compute-intensive study subjects used for evaluation are simulators or Simulink models that take between two to ten minutes to execute. When the execution times of the study subjects are in the order of a few minutes, it is still feasible to execute the SA-SBT and the baseline several times and draw conclusions using statistical tests. It is, however, important to note that, in these papers, the number of times experiments were repeated to account for randomness was rather low, i.e., between 10 to 30 times. This can negatively impact the accuracy of the statistical test results. Further, none of these papers performed any extensive meta-evaluations or hyper-parameter optimizations.

In the case of Menghi et. al. Menghi et al. (2020), a single simulation of the compute-intensive subject takes 1.5 hour, and the approach of the paper which involves the use of SI requires extensive hyper-parameter optimizations. Performing all the required experiments on their compute-intensive subject could take in the order of 50 years to complete. Hence, Menghi et. al. Menghi et al. (2020) perform a part of their experiments on non-compute-intensive subjects, measure the number of iterations required to detect failures and compute the *estimated execution time*, i.e., the estimated time that performing these numbers of iterations would require for compute-intensive models. The testing techniques are then compared by considering the *estimated execution time* instead of the same *execution time*.

**Metrics for assessing the performance of Pareto-based algorithms.** The quality indicators typically determine how fast the outputs of a Pareto-based algorithm converge towards an ideal or true Pareto front. These indicators may not directly relate to the main objective of a test generation algorithm though. Hence, there needs to be a justification as to why a quality indicator is used to assess a test generation algorithm. In our original study we selected three indicators, HV, Spread and GD, based on their widespread use and their collective ability to evaluate complementary aspects of convergence, spread, and uniformit Wu et al. (2022). As a result of our replication, we have reached at two observations related to the use of quality indicators for assessing Pareto-based test generation algorithms: (1) GD and Spread metrics are most meaningful for comparing population sets with equal sizes. Spread, while being normalised to the size of a population set, provides different values for two sets of different sizes, in which individuals are equally separated, which might lead to wrong interpretation of the spread for comparing populations of the same algorithm generated over subsequent runs. Even though GD does not have exactly the same problem, introduction of another individual to a population set of a small size might heavily worsen the metric value, and therefore, might also lead to incorrect interpretation of convergence over a run. We note that in contrast to classical Pareto optimization algorithms (e.g. NSGAII), NSGAII-DT allows for a flexible population size during the search. Relaxing the population size in search algorithms employed for testing is common not only for SA-SBT techniques but also in the many-objective search algorithms Abdessalem et al. (2018b); Panichella et al. (2015). The flexibility in the population size allows for more targeted exploration guided by surrogates, or to maintain a subset of the solutions on an optimal Pareto that are useful for the testing problem at hand. However, flexible population sizes can also affect metrics such as Spread and GD leading to potential misinterpretation of the results. (2) Among the three quality indicators used in our original study, the results related to HV in the replication study are the most faithful to those obtained from the original study. This might be partly due to the issue that Spread and GD are not suitable metrics for comparing Pareto-based algorithms with varying populations while HV is not dependent on the population size. Overall, in our experience, HV showed as a better metric for assessing our proposed Pareto-based SA-SBT test generation algorithms.

## 5.2. Surrogates for fault localization and repair

The main objective of existing SA-SBT approaches is to generate test cases that can effectively reveal failures. Empirical evidence from our subject papers Abdessalem et al. (2016, 2018a); Matinnejad et al. (2014); Menghi et al. (2020), the research reviewed in Section 3 and our two replication studies support the idea that the incorporation of surrogates into SBT significantly enhances the effectiveness and efficiency of automated testing in revealing failures. However, it is important to ensure that the techniques do not generate spurious failures, e.g., failures caused by invalid inputs and not due to system faults or failures that represent unrealistic situations. In our paper introducing NSGAII-DT, we utilized surrogate models to identify input conditions that led to failures, which we then validated against domain knowledge Abdessalem et al. (2018a). These conditions are valuable because they can help explain why failures occurred, which in turn can aid in root-cause analysis, fault localization, and bug fixing. However, we did not leverage these conditions to distinguish between valid and spurious failures, nor did we use them to identify root causes of failures. To the best of our knowledge, no prior research has explored the potential of surrogate models in explaining failure conditions. We believe that these models can serve as effective tools for identifying spurious failures, as well as developing fault localization and repair strategies.

## 5.3. Frameworks for large-scale experimentation

Further advancements in the area of SA-SBT require systematic evaluation frameworks that promote accuracy, transparency and reproducibility. Our first lessons-learned highlighted the challenge of defining proper evaluation metrics, but there is also the additional obstacle of limited platforms and benchmarks available for large-scale experimentation with SA-SBT methods. As discussed in our taxonomy (Section 2), SA-SBT methods may be applied in the MiL, SiL and HiL contexts, but the main focus of the research so far has been on MiL testing. The majority of MiL testing has been performed on case studies based on simulation-based testing of autonomous driving systems or Simulink models. A variety of simulators and ADAS components have been used in the literature for simulation-based testing of autonomous driving systems. The example simulators include: PreScan pre (2023), Pro-SiVIC Belbachir et al. (2012), CARLA Dosovitskiy et al. (2017), and BeamNG Gambi et al. (2019), and the ADAS components include both industrial and proprietary ADAS and open source DNN-enabled components.

A framework is needed to enable experimentation and evaluation of testing algorithms using different simulators and ADAS alternatives, as most current test generation algorithms are evaluated on a small subset, or even a single, simulator and ADAS. Our replication in Section 4.1 concerned reproducing results from an original study using the same simulator but two different ADAS. As shown there, a number of different factors need to be considered when we apply a heuristic algorithm to two different ADAS even when the simulator is the same. Other replication attempts in the context of ADAS testing include the work of Borg et. al. Borg et al. (2021) on replicating the results of a search-based testing algorithm on two different simulators (i.e., ProSiVic and PreScan) for the same ADAS component, and the work of Stocco et. al. Stocco et al. (2022) on generalizing the testing results obtained using an ADAS simulator to a physical platform. Both studies observe notable discrepancies between the testing results obtained from the different simulators and platforms. In particular, the latter work which is a rare example of transferring the testing result from MiL to HiL shows some major reality gaps between the virtual and physical world. A notable effort in evaluating and comparing MiL, SiL and HiL test setups is the recent work by Mandrioli et. al. Mandrioli et al.. This work which is performed in the context of drones demonstrates that the results obtained from MiL, SiL and HiL levels have notable differences, and in fact, these levels have complementary capabilities in revealing faults for CPS.

### 5.4. International Competitions

International competitions are a successful instrument to support experiment replication effectively. While participating in competitions requires significant effort, these competitions offer several benefits. We learned many lessons over our three-year participation in the ARCH competition that we summarize below.

*Forces tools updates.* Authors need to constantly update their tools to ensure specific formatting of their inputs and outputs to adhere to the rules of the competition. From our experience, we learned that this activity improved our tool. For example, participating in the competition in 2021 enabled us to extend the implementation of ARIsTEO to support property specifications that use locations. This activity extended the support of ARIsTEO to benchmarks AT51, AT52, AT53, and AT54.

*Supports experiment replication and verification of previous results.* Replicating previous experiments is not always a trivial and straightforward ac-

tivity. Replicating the experiments for AT51, AT52, AT54, and AT54 was not trivial. ARIsTEO is a plugin for S-TaLiRo. S-TaLiRo requires to specify these requirements using "control locations" (Pred(i).loc), i.e., specific constructs used to establish fitness metrics of hybrid systems. Despite S-TaLiRo being widely documented and the explanations being detailed and exhaustive, the documentation of this construct should be improved. Specifically, to analyze these requirements, we had to install MatlabBGL Mat (2022), a library containing a set of algorithms to work with graphs, and recompile the C++ files that are used to compute the fitness metric using an old version (1_36_0 boo (2022)) of Boost, a portable C++ and GCC library. However, this needed to be more precisely detailed since we had to reverse engineer the version of the library to be used: we could not use newer versions of Boost since they are not retro-compatible for this case. ARIsTEO participation was managed by different students across the years. We spent significant time in 2021 adding support for these requirements by extending the implementation of ARIsTEO and understanding how to install the different libraries. However, we did not precisely document our activity, and we had to perform the same reverse engineering process in 2022. We now have a detailed description of the libraries that need to be installed to support these requirements, add it to our documentation, and contact the developers of S-TaLiRo to extend their documentation. Finally, as detailed in Section 4.2, participating in consecutive years of the competition enables spotting configuration errors and improving the reliability of the published results.

*Discovering new results over time.* The set of benchmark models change over time: the ARCH competition is organized as part of the International Workshop on Applied Verification of Continuous and Hybrid Systems, which contains an explicit call for new benchmark models. For example, in 2022, the organizers removed the WT model from the set of benchmark models considered in the competition, and the pacemaker model Ayesh et al. (2022) was proposed at the workshop and is likely to be added to the set of benchmark models to be considered in future editions of the competition. The variation of the benchmark models enables a continuous comparison and assessment of the tools on models from different domains and with various characteristics (e.g., number of blocks).

*Supports networking and research collaborations.* The ARCH competition is a friendly competition with the primary purpose of keeping researchers and practitioners updated with the latest advancements in the field. It is designed to support and foster networking and research collaborations. This decision

ensures that participants are offered a friendly and cordial environment. We learned that this offered a great opportunity (especially for students) for networking, getting new research ideas, and establishing research collaborations.

Finally, we learned that organizing and managing the competition is a significant undertaking. The coordinator of the Falsification Category of the ARCH Competition (Gidon Ernst) did a significant job across the three years in which we participated to ensure the competition's success. Organizing the meetings to define the rules of the competition, collecting and analyzing the results, leading the writing of the competition report, and reporting the results at the workshop are a significant service to the community that should be acknowledged and appropriately rewarded.

Our taxonomy can support the organizers of the ARCH competition by supporting activities such as adding new benchmarks. For example, the organizers may extend the benchmarks to include models with more complexity and longer execution time. Further, the metrics we discussed earlier in this section can be used as part of the benchmark and competition since, currently, the number of iterations is the only metric used for the tool competition.

## 6. Conclusions

This paper reflects on four papers published between 2014 and 2021 Abdessalem et al. (2016, 2018a); Menghi et al. (2020); Matinnejad et al. (2014) that propose surrogate-assisted search-based testing (SA-SBT) techniques for autonomous systems. We developed a taxonomy based on our synthesis of different SA-SBT approaches. Our taxonomy identifies the main dimensions of SA-SBT solutions and we demonstrate how different SA-SBT solutions can be positioned along these dimensions. We report on two replication studies using an industrial advanced driver assistance system (ADAS) and a benchmark of Simulink models. The results of our taxonomy and replication studies highlight the need for improvement in the metrics used to evaluate SA-SBT, the potential for using surrogates in fault localization and repair, and the importance of benchmarking, international competitions and creating frameworks for large-scale experiments involving SA-SBT techniques.

## Acknowledgements

## References

, 2022. boost. `https://sourceforge.net/projects/boost/files/boost/1.36.0/`. Accessed: 2022-11-21.

, 2022. Design narma-l2 neural controller in simulink. `https://www.mathworks.com/help/deeplearning/ug/design-narma-l2-neural-controller-in-simulink.html`. Accessed: 2022-11-21.

, 2022. Matlabbgl. `https://www.mathworks.com/matlabcentral/fileexchange/10922-matlabbgl`. Accessed: 2022-11-21.

, 2023. Results replication NSGAII-DT. URL: `https://github.com/Leviathan321/reflection_study`.

, 2023. Simcenter Prescan. URL: `https://www.plm.automation.siemens.com/global/de/products/simcenter/prescan.html`. accessed: 2023-02-02.

Abdessalem, R.B., Nejati, S., Briand, L.C., Stifter, T., 2016. Testing advanced driver assistance systems using multi-objective search and neural networks, in: International Conference on Automated Software Engineering (ASE), IEEE/ACM. pp. 63–74.

Abdessalem, R.B., Nejati, S., Briand, L.C., Stifter, T., 2018a. Testing vision-based control systems using learnable evolutionary algorithms, in: Chaudron, M., Crnkovic, I., Chechik, M., Harman, M. (Eds.), International Conference on Software Engineering, (ICSE), ACM. pp. 1016–1026.

Abdessalem, R.B., Panichella, A., Nejati, S., Briand, L.C., Stifter, T., 2018b. Testing autonomous cars for feature interaction failures using many-objective search, in: Huchard, M., Kästner, C., Fraser, G. (Eds.), International Conference on Automated Software Engineering (ASE), ACM/IEEE. pp. 143–154.

AMS, M.M., Electromagnetism, P., 1989. System identification .

Annpureddy, Y., Liu, C., Fainekos, G., Sankaranarayanan, S., 2011. S-taliro: A tool for temporal logic falsification for hybrid systems, in: Tools and Algorithms for the Construction and Analysis of Systems, Springer. pp. 254–257.

ARCH, 2022 [Online]. International Competition on Verifying Continuous and Hybrid Systems. URL: https://cps-vo.org/group/ARCH/FriendlyCompetition.

Arrieta, A., Wang, S., Markiegi, U., Sagardui, G., Etxeberria, L., 2017. Search-based test case generation for cyber-physical systems, in: Congress on Evolutionary Computation (CEC), IEEE. pp. 688–697.

Ayesh, M., Mehan, N., Dhanraj, E., El-Rahwan, A., Opalka, S.E., Fan, T., Hamilton, A., Jacob, A.M., Sundarrajan, R.A., Widjaja, B., Menghi, C., 2022. Two simulink models with requirements for a simple controller of a pacemaker device, in: International Workshop on Applied Verification of Continuous and Hybrid Systems (ARCH22), EasyChair. pp. 18–25.

Beglerovic, H., Stolz, M., Horn, M., 2017. Testing of autonomous vehicles using surrogate models and stochastic optimization, in: International Conference on Intelligent Transportation Systems (ITSC), IEEE. pp. 1–6.

Belbachir, A., Smal, J.C., Blosseville, J.M., Gruyer, D., 2012. Simulation-Driven Validation of Advanced Driving-Assistance Systems. Procedia - Social and Behavioral Sciences 48, 1205–1214. doi:10.1016/j.sbspro.2012.06.1096.

Bittanti, S., 2019. Model identification and data analysis. John Wiley & Sons.

Borg, M., Abdessalem, R.B., Nejati, S., Jegeden, F., Shin, D., 2021. Digital twins are not monozygotic - cross-replicating ADAS testing in two

industry-grade automotive simulators, in: Conference on Software Testing, Verification and Validation (ICST), IEEE. pp. 383–393.

Chen, T., Li, M., Yao, X., 2020. How to evaluate solutions in pareto-based search-based software engineering? A critical review and methodological guidance. CoRR abs/2002.09040. URL: `https://arxiv.org/abs/2002.09040`, `arXiv:2002.09040`.

Clarke, E.M., Biere, A., Raimi, R., Zhu, Y., 2001. Bounded model checking using satisfiability solving. Formal Methods in System Design 19, 7–34.

Deb, K., Pratap, A., Agarwal, S., Meyarivan, T., 2002. A fast and elitist multiobjective genetic algorithm: NSGA-II. IEEE Transactions on Evolutionary Computation 6, 182–197. doi:`10.1109/4235.996017`.

Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., Koltun, V., 2017. Carla: An open urban driving simulator. arXiv preprint arXiv:1711.03938 .

Dybå, T., Maiden, N.A.M., Glass, R., 2014. The reflective software engineer: Reflective practice. IEEE Softw. 31, 32–36.

Ernst, G., Arcaini, P., Bennani, I., Donzé, A., Fainekos, G., Frehse, G., Mathesen, L., Menghi, C., Pedrielli, G., Pouzet, M., Yaghoubi, S., Yamagata, Y., Zhang, Z., 2020. ARCH-COMP 2020 category report: Falsification, in: International Workshop on Applied Verification of Continuous And Hybrid Systems, EasyChair. pp. 140–152.

Ernst, G., Arcaini, P., Fainekos, G., Formica, F., Inoue, J., Khandait, T., Mahboob, M.M., Menghi, C., Pedrielli, G., Waga, M., Yamagata, Y., Zhang, Z., 2022. Arch-comp 2022 category report: Falsification with ubounded resources, in: International Workshop on Applied Verification of Continuous and Hybrid Systems (ARCH22), EasyChair. pp. 204–221. doi:`10.29007/fhnk`.

Ernst, G., Sedwards, S., Zhang, Z., Hasuo, I., 2018. Fast falsification of hybrid systems using probabilistically adaptive input. arXiv preprint arXiv:1812.04159 .

G. Ernst, P. Arcaini, I.B.A.C., Donzé, A., Fainekos, G., Frehse, G., Gaaloul, K., Inoue, J., Khandait, T., Mathesen, L., Menghi, C., Pedrielli, G., Pouzet, M., Waga, M., Yaghoubi, S., Yamagata, Y., Zhang, Z., 2021.

ARCH-COMP category report: Falsification with validation of results, in: Workshop on Applied Verification of Continuous and Hybrid Systems, EasyChair. pp. 133–152.

Gambi, A., Müller, M., Fraser, G., 2019. AsFault: testing self-driving car software using search-based procedural content generation, in: Atlee, J.M., Bultan, T., Whittle, J. (Eds.), International Conference on Software Engineering: Companion Proceedings (ICSE), IEEE / ACM. pp. 27–30.

Haq, F.U., Shin, D., Briand, L., 2022. Efficient online testing for dnn-enabled systems using surrogate-assisted and many-objective optimization, in: International Conference on Software Engineering (ICSE 2022), pp. 811–822.

Heidlauf, P., Collins, A., Bolender, M., Bak, S., 2018. Verification challenges in F-16 ground collision avoidance and other automated maneuvers, in: Frehse, G. (Ed.), ARCH18. 5th International Workshop on Applied Verification of Continuous and Hybrid Systems, EasyChair. pp. 208–217. doi:`10.29007/91x9`.

Henard, C., Papadakis, M., Perrouin, G., Klein, J., Traon, Y.L., 2013. PLEDGE: a product line editor and test generation tool, in: International Software Product Line Conference co-located workshops (SPLC'13), ACM. pp. 126–129.

Hoxha, B., Abbas, H., Fainekos, G., 2015. Benchmarks for temporal logic requirements for automotive systems, in: Frehse, G., Althoff, M. (Eds.), ARCH14-15. 1st and 2nd International Workshop on Applied veRification for Continuous and Hybrid Systems, EasyChair. pp. 25–30. doi:`10.29007/xwrs`.

Hu, J., Lygeros, J., Sastry, S., 2000. Towards a theory of stochastic hybrid systems, in: International Workshop on Hybrid Systems: Computation and Control, Springer. pp. 160–173.

Humeniuk, D., Antoniol, G., Khomh, F., 2021. Data driven testing of cyber physical systems. CoRR abs/2102.11491. URL: `https://arxiv.org/abs/2102.11491`, `arXiv:2102.11491`.

Humeniuk, D., Khomh, F., Antoniol, G., 2022. A search-based framework for automatic generation of testing environments for cyber-physical systems. Inf. Softw. Technol. 149, 106936. doi:`10.1016/j.infsof.2022.106936`.

Innes, C., Ramamoorthy, S., 2022. Automated testing with temporal logic specifications for robotic controllers using adaptive experiment design, in: International Conference on Robotics and Automation (ICRA), IEEE. pp. 6814–6821. doi:`10.1109/ICRA46639.2022.9811579`.

Jin, X., Deshmukh, J.V., Kapinski, J., Ueda, K., Butts, K., 2014. Powertrain control verification benchmark, in: International Conference on Hybrid Systems: Computation and Control, ACM. pp. 253–262. doi:`10.1145/2562059.2562140`.

Jin, Y., 2011. Surrogate-assisted evolutionary computation: Recent advances and future challenges. Swarm and Evolutionary Computation 1, 61–70.

Luke, S., 2013. Essentials of Metaheuristics. second ed., Lulu. Available for free at http://cs.gmu.edu/~sean/book/metaheuristics/.

Mandrioli, C., Carlsson, M.N., Maggio, M., . Testing abstractions for cyber-physical control systems. Submitted for publication.

Matinnejad, R., Nejati, S., Briand, L.C., Bruckmann, T., 2014. Mil testing of highly configurable continuous controllers: scalable search using surrogate models, in: Crnkovic, I., Chechik, M., Grünbacher, P. (Eds.), International Conference on Automated Software Engineering (ASE), ACM/IEEE. pp. 163–174.

McKay, M.D., Beckman, R.J., Conover, W.J., 1979. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. Technometrics 21, 239–245. URL: `http://www.jstor.org/stable/1268522`.

Menghi, C., Nejati, S., Briand, L., Parache, Y.I., 2020. Approximation-refinement testing of compute-intensive cyber-physical models: An approach based on system identification, in: IEEE/ACM 42nd International Conference on Software Engineering (ICSE), pp. 372–384.

Panichella, A., Kifetew, F.M., Tonella, P., 2015. Reformulating branch coverage as a many-objective optimization problem, in: International Conference on Software Testing, Verification and Validation (ICST), IEEE Computer Society. pp. 1–10.

Pedrielli, G., Khandait, T., Chotaliya, S., Thibeault, Q., Huang, H., Castillo-Effen, M., Fainekos, G., 2021. Part-x: A family of stochastic algorithms for search-based test generation with probabilistic guarantees. CoRR abs/2110.10729. URL: `https://arxiv.org/abs/2110.10729`, `arXiv:2110.10729`.

Schuler, S., Adegas, F.D., Anta, A., 2017. Hybrid modelling of a wind turbine, in: Frehse, G., Althoff, M. (Eds.), ARCH16. International Workshop on Applied Verification for Continuous and Hybrid Systems, EasyChair. pp. 18–26. doi:`10.29007/tf1p`.

Stocco, A., Pulfer, B., Tonella, P., 2022. Mind the gap! a study on the transferability of virtual vs physical-world testing of autonomous driving systems. IEEE Transactions on Software Engineering , 1–13doi:`10.1109/tse.2022.3202311`.

Waga, M., 2020. Falsification of cyber-physical systems with robustness-guided black-box checking, in: International Conference on Hybrid Systems: Computation and Control, ACM. pp. 11:1–11:13.

Wang, Y., Yu, R., Qiu, S., Sun, J., Farah, H., 2022. Safety performance boundary identification of highly automated vehicles: A surrogate model-based gradient descent searching approach. IEEE Transactions on Intelligent Transportation Systems , 1–12doi:`10.1109/TITS.2022.3191088`.

Wu, J., Arcaini, P., Yue, T., Ali, S., Zhang, H., 2022. On the preferences of quality indicators for multi-objective search algorithms in search-based software engineering. Empirical Software Engineering 27, 144. doi:`10.1007/s10664-022-10127-4`.

Yaghoubi, S., Fainekos, G., 2019. Gray-box adversarial testing for control systems with machine learning components, in: International Conference on Hybrid Systems: Computation and Control (HSCC).

Yamagata, Y., Liu, S., Akazaki, T., Duan, Y., Hao, J., 2021. Falsification of cyber-physical systems using deep reinforcement learning. IEEE Transactions on Software Engineering 47, 2823–2840. doi:`10.1109/TSE.2020.2969178`.

Zeller, A., 2017. Search-based testing and system testing: A marriage in heaven, in: International Workshop on Search-Based Software Testing, (SBST@ICSE), IEEE/ACM. pp. 49–50.

Zhang, Z., Arcaini, P., 2021. Gaussian process-based confidence estimation for hybrid system falsification, in: Huisman, M., Pasareanu, C.S., Zhan, N. (Eds.), Formal Methods (FM), Springer. pp. 330–348.

Zhang, Z., Lyu, D., Arcaini, P., Ma, L., Hasuo, I., Zhao, J., 2021. Effective hybrid system falsification using monte carlo tree search guided by QB-robustness, in: Silva, A., Leino, K.R.M. (Eds.), Computer Aided Verification, Springer. pp. 595–618.

Zhong, Z., Kaiser, G.E., Ray, B., 2021. Neural network guided evolutionary fuzzing for finding traffic violations of autonomous vehicles. CoRR abs/2109.06126. URL: https://arxiv.org/abs/2109.06126, arXiv:2109.06126.