# Test Generation Strategies for Building Failure Models and Explaining Spurious Failures

BAHARIN A. JODAT, University of Ottawa, Canada
ABHISHEK CHANDAR, University of Ottawa, Canada
SHIVA NEJATI, University of Ottawa, Canada
MEHRDAD SABETZADEH, University of Ottawa, Canada

Test inputs fail not only when the system under test is faulty but also when the inputs are invalid or unrealistic. Failures resulting from invalid or unrealistic test inputs are spurious. Avoiding spurious failures improves the effectiveness of testing in exercising the main functions of a system, particularly for compute-intensive (CI) systems where a single test execution takes significant time. In this paper, we propose to build failure models for inferring interpretable rules on test inputs that cause spurious failures. We examine two alternative strategies for building failure models: (1) machine learning (ML)-guided test generation and (2) surrogate-assisted test generation. *ML-guided test generation* infers boundary regions that separate passing and failing test inputs and samples test inputs from those regions. *Surrogate-assisted test generation* relies on surrogate models to predict labels for test inputs instead of exercising all the inputs. We propose a novel surrogate-assisted algorithm that uses multiple surrogate models simultaneously, and dynamically selects the prediction from the most accurate model. We empirically evaluate the accuracy of failure models inferred based on surrogate-assisted and ML-guided test generation algorithms. Using case studies from the domains of cyber-physical systems and networks, we show that our proposed surrogate-assisted approach generates failure models with an average accuracy of 83%, significantly outperforming ML-guided test generation and two baselines. Further, our approach learns failure-inducing rules that identify genuine spurious failures as validated against domain knowledge.

CCS Concepts: • **Software and its engineering** → **Search-based software engineering**; **Empirical software validation**; • **Computing methodologies** → **Machine learning algorithms**.

Additional Key Words and Phrases: Search-based testing, Machine learning, Surrogate models, Failure models, Test-input validity, and Spurious failures.

## 1 INTRODUCTION

Traditionally, software testing has been concerned with finding inputs that reveal failures in the system under test (SUT). However, failures do not always indicate faults in the SUT. Instead, failures may arise because the test inputs are *invalid* or *unrealistic*. For example, in an airplane autopilot system, a failure indicating that the ascent requirement fails for a test input that points the plane

Authors' addresses: Baharin A. Jodat, balia034@uottawa.ca, University of Ottawa, Canada; Abhishek Chandar, achan260@uottawa.ca, University of Ottawa, Canada; Shiva Nejati, snejati@uottawa.ca, University of Ottawa, Canada; Mehrdad Sabetzadeh, m.sabetzadeh@uottawa.ca, University of Ottawa, Canada.

nose downward would be invalid. This is because the failure is caused by an unmet environment assumption: the nose should be upward for ascent. For another example, consider a network-management system. In such a system, quality-of-service requirements will inevitably fail for unrealistic test inputs that overwhelm the network beyond its capacity. Environment assumptions capture the expected conditions for a system's operational environment [49]. Attempting to test a system for a more general environment than its expected operational environment may lead to overly pessimistic testing and verification results. We refer to failures arising from test inputs violating the system's environment assumptions as *spurious failures*. It is often the case that environment assumptions are not fully known for software systems [83]; therefore, it is difficult to determine whether a failure is indeed spurious.

Automated random testing (fuzzing) [75] becomes more effective in exercising the main functions of a system if the fuzzer avoids spurious failures [51]. Spurious failures particularly pose a challenge for compute-intensive (CI) systems, where a single test execution takes significant time to complete. For CI systems, we want to use the limited testing time budget to generate valid inputs that exercise the system's main functions. A promising approach for identifying spurious failures is to build *failure models* [65–67]. Failure models provide conditions that explain the circumstances of failures and describe when a failure occurs and when it does not [65]. Failure models can infer rules leading to and only to failures. These rules are candidates to be validated against domain knowledge to determine whether the failures that the rules identify are spurious.

Recent research on synthesizing input grammars [32, 33, 67, 68, 89] and abstracting failure-inducing inputs [51, 65, 66] aims to understand the circumstances of different failures. These approaches start from an example failure and iteratively generate more tests to learn the input conditions that lead to that failure. The tests are generated via fuzz testing with or without an input grammar. These approaches are geared towards systems with string inputs, where oracles are typically binary (pass/fail) verdicts. However, these approaches are not optimized for systems with numeric inputs, where the inputs are not governed by grammars and where quantitative fitness functions, developed based on system requirements, are used to determine the degree to which test inputs pass or fail. These quantitative fitness functions enable exploration of input space using multiple test-generation heuristics and learning algorithms, resulting in test sets with sufficient information to infer candidate rules for identifying spurious failures.

This paper proposes a framework to infer failure models for compute-intensive (CI) systems with numeric inputs. Examples of such systems include cyber-physical systems (CPS) and network systems. We follow a data-driven approach and infer failure models by harvesting information from a set of test inputs. To generate such sets, one can use either *explorative* or *exploitative* search methods [70]. The former attempts to sample the entire search space, whereas the latter attempts to sample the most informative regions of the search space. The challenge with the explorative approach is that we need to collect and execute many test inputs from the search space to determine if they pass or fail. For CI systems, this takes significant time and may become infeasible. The challenge with the exploitative approach is that one needs effective guidance for sampling within large and multi-dimensional search spaces.

Machine learning (ML) has been used for improving the effectiveness and efficiency of both explorative and exploitative search [35, 47, 63, 69, 71, 73, 86]. For explorative search, *surrogate-assisted test generation* relies on ML to predict verdicts for test inputs instead of executing them all [35, 62, 71, 78]. Using a quantitative surrogate, one can forego system executions when the predicted verdicts remain valid after offsetting prediction errors. Otherwise, we execute the SUT and use the results from the executed test inputs to refine the surrogate. As for exploitative search, *ML-guided test generation* aims to infer boundary regions that separate passing and failing test inputs and to subsequently sample test inputs from those regions [47, 69]. The intuition is that

tests sampled in the boundary regions are more informative for identifying failures and can be used to further refine the boundary regions. Both approaches provide a set of labelled test inputs from which one can infer failure models using techniques such as decision-rule learning [91]. Human experts must nonetheless review and validate the resulting rules to determine whether they represent genuine spuriousness. The use of interpretable ML techniques such as decision-rule learning allows failure models to be expressed as easily understandable rules linked to system inputs, making them ideal for human interpretation.

While surrogate-assisted and ML-guided test generation algorithms have been previously used to generate individual test inputs [35, 78] , their efficacy in generating failure models remains unexplored. Specifically, earlier work strands [35, 78] employ machine learning to more effectively steer test generation towards areas within the search space that are likely to contain the most severe failures. In this paper, we use machine learning to devise new test generation algorithms, with the aim of inferring failure models for systems that have numeric inputs.We evaluate the resulting failure models against those produced by baselines. Our evaluation answers two main questions: (1) How accurate are the failure models generated by the surrogate-assisted and ML-guided techniques in predicting failures? (2) How useful are failure models for identifying spurious failures? We use two kinds of study subjects in our evaluation: (i) A benchmark of four CPS Simulink subjects with 12 requirements that are non-compute intensive (non-CI). (ii) Two industrial CI systems, one from the CPS and the other from the network domain. In summary, we make the following contributions:

**(1)** We propose a data-driven framework for inferring failure models for systems with numeric inputs including CPS and network systems (Section 3).

**(2)** We propose a dynamic surrogate-assisted algorithm that uses multiple surrogate models simultaneously during search, and dynamically selects the prediction from the most accurate model (Section 3.2). Our evaluation performed based on seven surrogate-model types in the literature [43, 44, 78, 86] shows that, compared to using surrogate models individually, our dynamic surrogate-assisted algorithm provides the best trade-off between accuracy and efficiency by generating datasets that are at least 33% larger while being at least 28% more accurate (RQ1 in Section 4.2).

**(3)** We compare the accuracy of failure models obtained using our dynamic surrogate-assisted approach against two ML-guided techniques as well as two baselines. One baseline is random-search, and the other is an adaptation of a state-of-the-art approach that generates failure models for systems with structured inputs [65]. Our results show that our dynamic surrogate-assisted algorithm yields failure models with an average accuracy, precision, and recall of 83%, 72%, and 88%, respectively, significantly outperforming the ML-guided algorithms and the baselines (RQ2 in Section 4.3 and RQ3 in Section 4.4).

**(4)** We demonstrate that failure models built using our dynamic surrogate-assisted algorithm generate useful rules for identifying spurious failures in our CI subjects, as validated by domain knowledge (RQ4 in Section 4.5).

**(5)** We present lessons learned based on our findings: The first lesson summarizes the advantages of using decision rules for building failure models. The second lesson highlights the limitations of focusing testing on finding individual failures and why failure models provide better insights about the effectiveness of testing algorithms.

It is essential for systems to handle all potential inputs including those that violate environment assumptions, and hence, are invalid. Spurious failures caused by invalid test inputs indicate a need for additional safeguards against invalid inputs that may be generated, among other sources, by human operator errors or malfunctioning hardware components, such as inaccurate sensor data. However, these invalid test inputs do not exercise the core functionality of a system. While ensuring that a given system is safeguarded against invalid inputs is crucial, the inability to identify spurious

failures can distort our understanding of the system's capabilities. This may also lead to misplaced confidence in a testing strategy that reveals numerous failures, yet offers little insight into the system's primary functions.

**Organization.** Section 2 motivates the need for identifying spurious failures. Section 3 presents our data-driven framework for inferring failure models and presents alternative surrogate-assisted and ML-guide algorithms for building failure models. Section 4 presents an evaluation of these algorithms. Section 5 outlines the main lessons learned from the research. Section 6 compares with related work. Section 7 summarizes the paper and suggests directions for future work.

## 2   MOTIVATION

Using two real-world, compute-intensive (CI) systems, we motivate the need for identifying spurious failures. These systems, both of which are open-source, are a Network Traffic Shaping System (NTSS) [57, 64] and an autopilot system [1].

NTSS is typically deployed on routers to ensure high network performance for real-time streaming applications such as teleconferencing (e.g., Zoom). Without NTSS, voice and video packets may be transmitted out of sequence or with delays. As a result, users may experience choppy or freezing voice/video. Due to the increasing remote-working practices, multiple streaming applications may be running at the same time in homes and small-office settings. This has made systematic testing of NTSS essential as a way to ensure that networks meet their quality-of-experience requirements.

NTSS works by dividing the total network bandwidth into classes with different priorities. The higher-priority classes are typically used for transmitting time-sensitive, streaming voice and video. To test the performance of an NTSS, we assign data flows with different bandwidth values to different NTSS classes. The purpose is to ensure that NTSS is configured optimally and can maintain good performance even when a high volume of traffic flows through its different classes. When we stress-test an NTSS, no matter how well-designed the NTSS is, we expect the quality of experience to deteriorate and become unacceptable eventually. Test inputs that stress NTSS beyond a certain limit deterministically fail and do not help reveal flaws or suboptimality in the NTSS design. Our approach in this paper infers the limit on the traffic that can flow through different NTSS classes without compromising the quality of experience. For an NTSS setup with eight classes from class0 to class7, we learn the following rule specifying failing test inputs:

> r1: IF (class5+ class6+ class7 > 0.75 · threshold) THEN FAIL

In the above rule, threshold is the sum of the maximum bandwidths of classes 5, 6, and 7. As we discuss in Section 4.5, we validate Rule r1 with a domain expert and confirm that failures specified by this rule are indeed spurious. Rule r1 indicates that attempting to simultaneously utilize classes 5, 6 and 7 more than 75% of their maximum ranges would compromise quality of experience for the entire network. Rule r1 helps domain experts in at least two ways: (1) it informs them that test inputs that satisfy the rule are spurious, since such test cases do not reveal design faults, and (2) it provides experts with data-driven evidence that the cumulative utilization of classes 5, 6, and 7 should be kept below the identified limit.

Our second case study is an autopilot system from Lockheed Martin's benchmark of challenge Simulink models [38]. This autopilot system is expected to satisfy the following requirement: $\varphi$= "*When the autopilot is enabled, the aircraft should reach the desired altitude within 500 seconds in calm air*". When we test the autopilot by fuzzing, we find several test inputs that violate this requirement and several test inputs that satisfy it. It is however unclear whether the failures are due to faults in the system or due to missing or unknown assumptions on the system inputs. Failures

caused by missing or unknown assumptions would be spurious. As we discuss in Section 4.5, we identify the following rule as one that indicates spurious failures:

> r2: IF (PitchWheel[0..300] ≤ -28 ∧ Throttle[0..300] ≤ 0.1) THEN FAIL

Here, Throttle[0..300] is the boost applied to the engine by the pilot during the first 300s, and PitchWheel[0..300] is the upward or downward degree of the aircraft nose, again during the first 300s. Note that both Throttle and PitchWheel are signals over time. In order to validate r2, we examined the handbook of the De Havilland Beaver aircraft [31]. According to the handbook, for this aircraft type, to satisfy requirement $\varphi$, the pilot should manually adjust the throttle boost (Throttle) to a sufficiently high value. The handbook further states that to be able to ascend, the plane's nose should not be pointing downward. That is, r2 describes a situation where the pre-conditions for $\varphi$ are not met. Hence, the tests in these ranges are expected to fail and are uninteresting for revealing system faults. Further, r2 can be used for implementing safeguards against misuse by the human operator (pilot).

## 3 GENERATING FAILURE MODELS

Figure 1 shows our framework for generating failure models. The inputs to our framework are: (1) an executable system or simulator $S$, (2) the input-space representation $\mathcal{R}$ for $S$, and (3) a quantitative fitness function $F$ for each requirement of $S$; an example requirement, $\varphi$, for autopilot was given in Section 2. The full set of requirements for our case studies is available in our supplementary material [6]. We make the following assumptions about the search input space ($\mathcal{R}$) and the fitness function ($F$):

- **A1** We assume that the system inputs are variables of type real or enumerate. For each real variable, the range of the values that the variable can take is bounded by an upper bound and a lower bound.
- **A2** For each requirement of $S$, we have a fitness function $F$ based on which a pass/fail verdict can be derived for any test input. Further, the value of $F$ differentiates among the pass test inputs, those that are more acceptable, and among the fail test inputs, those that trigger more severe failures. Specifically, we assume that the range of $F$ is an interval $[-a, b]$ of $\mathbb{R}$. For a given test, $F(t) \geq 0$ iff $t$ is passing, and otherwise, $t$ is failing. The closer $F(t)$ is to $b$, the higher the confidence that $t$ passes; and the closer $F(t)$ is to $-a$, the higher the confidence that $t$ fails.

Assumptions **A1** and **A2** are common for CPS models expressed in Simulink [29, 38, 74, 87], automated driving systems [78], and network-management systems [64], and are valid for all the case studies we use in our evaluation.

As discussed in Section 1, we examine two alternative test-generation approaches for building failure models: *surrogate-assisted* and *ML-guided*. Both approaches can be captured using the framework shown in Figure 1: The preprocessing phase generates a set of test inputs labelled with fitness values. The main loop takes the test-input set created by the preprocessing phase, and trains a model. When the framework is instantiated for surrogate-assisted test generation, the model predicts fitness values for the generated test inputs. When the framework is instantiated for ML-guided test generation, the model guides test-input sampling. The main loop extends the test-input set using the trained model while also refining the model based on newly generated tests. After the main loop terminates, the framework uses the test-input set to train, using decision-rule learners, a failure model. In the remainder of this section, we detail each step of the framework shown in Figure 1.
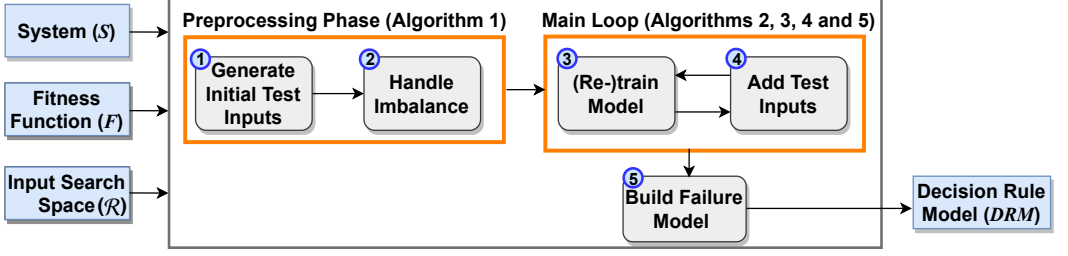
Fig. 1. Our framework for generating failure models. The main loop of the framework, i.e., steps 3 and 4, can be realized using two alternative test-generation strategies: (1) Surrogate-assisted test generation, or (2) ML-guided test generation. For surrogate-assisted test generation, one can use either Algorithm 2, which is based on an individual surrogate model, or Algorithm 3, which is based on our proposed dynamic model. For ML-guided test generation, one can use either Algorithm 4, which uses regression trees to identify the boundary regions between passing and failing test cases, or Algorithm 5, which uses logistic regression for the same purpose.

---

**Algorithm 1** The preprocessing phase of Figure 1

---

**Input** $S$: System
**Input** $\mathcal{R} = \{R_1, \ldots, R_n\}$: Ranges for input variables $v_1$ to $v_n$
**Input** $F$: Fitness Function
**Input** $d$: The budgeted dataset size
**Output** $DS$: A set of test inputs and their fitness values

1: $DS^i \leftarrow$ GenerateTests($\mathcal{R}$, $d/2$); //(Adaptive) Random Testing
2: $DS^l \leftarrow$ Execute($S$, $DS^i$, $F$); //Compute fitness values
3: $DS^b \leftarrow$ HandleImbalance($DS^l$); //Use SMOTE to generate synthetic test inputs
4: $DS^{b,l} \leftarrow$ Execute($S$, $DS^b$, $F$); //Compute fitness values of the tests generated by SMOTE
5: $DS \leftarrow DS^{b,l} \cup DS^l$; //Combine the test inputs generated by adaptive random testing and SMOTE
   along with their fitness values to form a dataset
6: **return** $DS$

---

## 3.1 Preprocessing Phase

Algorithm 1 describes the preprocessing phase that generates the initial dataset for training a model to be used in the main loop of Figure 1. This algorithm first randomly generates half ($d/2$) of the budgeted test inputs and computes a fitness value for each test input by executing the test input using $S$ (lines 1-2). Since ML models perform poorly when the training set is imbalanced, we attempt to address any potential imbalance before using the data for ML training [91] (line 3). In our work, the imbalance, if one exists, is between the pass and the fail classes. We use the well-known synthetic minority over-sampling technique (SMOTE) for addressing imbalance [39]. Let *minor* (resp. *major*) be the number of tests in the minority (resp. majority) class. SMOTE over-samples the minority class by taking each minority-class sample and introducing synthetic examples along the line segments joining any/all of the $k$ minority-class nearest neighbours [39]. The process is repeated until we have $m = major - minor$ new such tests.

We discard the labels from SMOTE and instead execute the tests to compute their actual fitness values (line 4). We discard the labels provided by SMOTE, since SMOTE categorizes test inputs as pass or fail. Instead, we require test inputs to be labelled with their quantitative fitness values; this enables us to train regression ML models in Step 3 of our approach in Figure 1. The final

dataset ($DS$) is returned at the end (line 6). Although not shown in Algorithm 1, to have exactly $d$ tests in $DS$, we generate the remaining ($d/2 - m$) tests randomly. Generating these remaining tests randomly does not introduce a new imbalance problem because if random test generation leads to major imbalance, then $m$ is already close to $d/2$ and only a few additional tests need to be generated. Otherwise, a small $m$ indicates that random testing is relatively balanced; in that case, no special provision is necessary for imbalance mitigation in the randomly generated datatset. Since we discard the labels generated by SMOTE and compute the actual labels, the imbalance problem may in principle persist even after applying SMOTE. For our experiments (Section 4), most synthetic samples generated by SMOTE indeed belong to the minority class. Our preprocessing therefore successfully addresses imbalance in our case studies.

## 3.2 Main Loop

The main data-generation loop is realized via two alternative algorithms, described below: surrogate-assisted and ML-guided.

The goal of this approach is to use surrogates to predict fitness values for some test inputs and thus not execute the system for all test inputs. Hence, surrogates help explore a larger portion of the input space and generate larger test sets.

*3.2.1 Surrogate-Assisted Test Generation.* Figure 2 illustrates the surrogate-assisted test generation process, which takes the same inputs as the framework in Figure 1. The output is a labelled dataset, $DS$, used in Step 5 of Figure 1 to build failure models. The procedure in Figure 2 is as follows:

(1) Start with preprocessing (Algorithm 1) to produce an initial dataset (Step 1 of Figure 2).
(2) Train a surrogate model with the initial dataset (Step 2 of Figure 2).
(3) Generate a new test input (Step 3 of Figure 2) and predict its fitness value using the surrogate model (Step 4 of Figure 2).
(4) Calculate a confidence interval for the predicted fitness value (Step 5 of Figure 2).
(5) If the prediction is not accurate, run system $S$ to obtain the actual fitness value (Step 7 of Figure 2), add the test input along with its actual fitness value to the dataset (Step 8 of Figure 2), and return to Step 2 of Figure 2 to re-train the surrogate model.
(6) If the prediction is accurate; add the test input along with its predicted fitness value to the dataset (Step 6 of Figure 2); and, return to Step 3 of Figure 2.

Note that if we execute system $S$ for a test input, as mentioned in item (5) above, the surrogate model is retrained using the dataset updated with this new test execution. In contrast, if system execution is not required, as in item (6) above, retraining the surrogate model is unnecessary.

To demonstrate the calculation of the confidence interval described in item (4) above, consider the example provided in Figure 3. In this figure, two sample test inputs, denoted as $t$ and $t'$, are shown. Suppose the predicted fitness values are $\bar{F}(t) = 8$ and $\bar{F}(t') = -1$, indicating a pass label for $t$ and a fail label for $t'$, respectively. Assume that the prediction error, $e$, of the surrogate model is 2. The confidence interval is calculated as $\bar{F} \pm e$. That is, the confidence interval for $\bar{F}(t)$ is $[6, 10]$, while the confidence interval for $\bar{F}(t')$ is $[-3, 1]$. Using the confidence intervals, we determine whether to execute system $S$ for $t$ and $t'$. For input $t$, the confidence interval of $\bar{F}(t)$ falls entirely within the positive range. Hence, even after accounting for error, we still label the test input $t$ as a pass. Therefore, there is no need for system $S$ to be executed for input $t$, since $t$ can be confidently labelled as pass (item (6) above). However, for input $t'$ the confidence interval of $\bar{F}(t')$ spans both positive and negative values. This indicates that a label cannot be confidently assigned to $t'$. As a result, system $S$ needs to be executed for $t'$ to obtain its actual fitness value and an accurate verdict (item (5)).
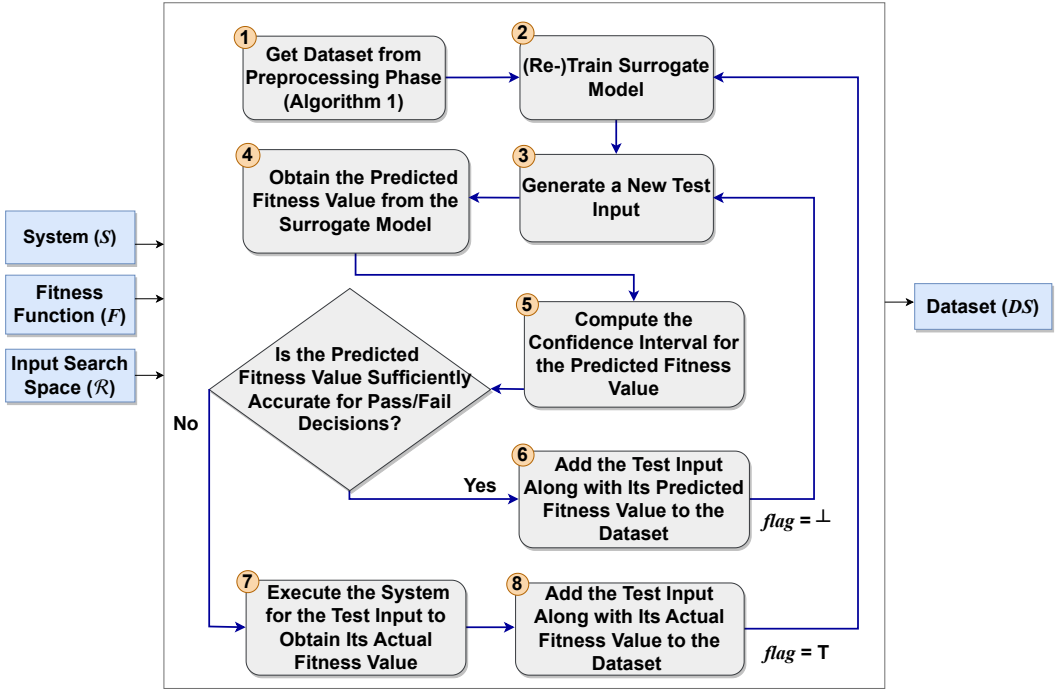
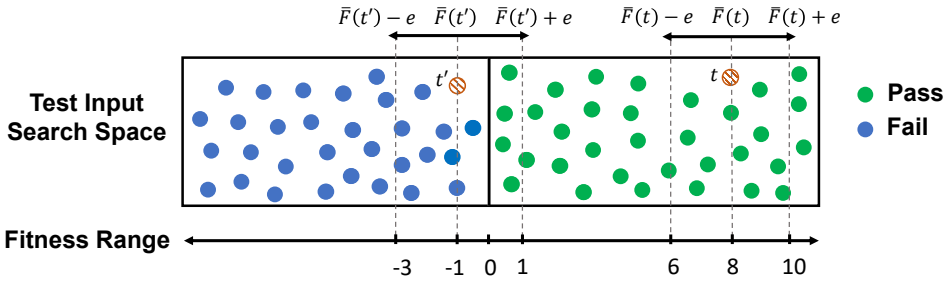Fig. 2. Illustration of the work flow of Algorithm 2



Fig. 3. Illustration of how to calculate the confidence interval for predicted fitness values to determine whether a predicted fitness value is accurate. Specifically, the figure shows confidence intervals, i.e., $\bar{F} \pm e$, of the predicted fitness values for two test inputs $t$ and $t'$. For the test input $t$ where the predicted fitness confidence interval remains entirely in the positive range (i.e., the verdict of the test inputs in this range is pass), we do not execute the system since we can say that the test input is passing even after accounting for error $e$. However, for the test input $t'$, the predicted fitness confidence interval spans both positive and negative values (i.e., it includes test inputs with both pass and fail verdicts). Therefore, we need to execute the system for $t'$ to obtain its actual fitness value.

The pseudo code of the surrogate-assisted test generation is provided in Algorithm 2. Below, we discuss each line of this algorithm in detail.

Line 1 of Algorithm 2 uses Algorithm 1 to obtain an initial dataset, *DS*:

---

**Algorithm 2** Test generation using surrogates

---

**Input** $S$: System
**Input** $\mathcal{R} = \{R_1, \ldots, R_n\}$: Ranges for input variables $v_1$ to $v_n$
**Input** $F$: Fitness Function
**Input** INITIALDATASETSIZE: size of initial dataset
**Output** $DS$: A dataset to train failure models

1: $DS \leftarrow$ Preprocessing($S$, $\mathcal{R}$, $F$, INITIALDATASETSIZE); //Run Algorithm 1 to generate an initial dataset
2: $DS^l \leftarrow DS$; $flag \leftarrow \top$;
3: **while** (execution budget remains) **do**
4:   **if** ($flag$)
5:    ($SM, e$) $\leftarrow$ Train($DS^l$); //Training $SM$; $e$ is the error
6:   **end**
7:   $t \leftarrow$ GenerateTests($\mathcal{R}$, 1); //Generate one test input
8:   $\bar{F}(t) \leftarrow SM(t)$; // Use $SM$ to predict a fitness value for $t$
9:   **if** ($\exists \sim \in \{\geq, <\} \cdot (\bar{F}(t) \sim 0) \wedge (\bar{F}(t) \pm e \sim 0)$ //Calculate the confidence interval of the predicted fitness value to decide whether system $S$ needs to be executed for $t$ or not
10:    $DS \leftarrow DS \cup \{\langle t, \bar{F}(t)\rangle\}$; //Add the test input along with its predicted fitness value to the dataset
11:    $flag \leftarrow \bot$; //No SM re-training is needed
12:   **else**
13:    $\{\langle t, F(t)\rangle\} \leftarrow$ Execute($S$, $\{t\}$, $F$); //Obtain the actual fitness value of $t$ by executing system $S$
14:    $DS^l \leftarrow DS^l \cup \{\langle t, F(t)\rangle\}$; //Add the test input along with its actual fitness value to the dataset
15:    $DS \leftarrow DS \cup \{\langle t, F(t)\rangle\}$; $flag \leftarrow \top$; //SM re-training is needed
16:   **end**
17: **end**
18: **return** $DS$

---

1: $DS \leftarrow$ Preprocessing($S$, $\mathcal{R}$, $F$, INITIALDATASETSIZE);

On line 2, $DS^l$ is created to maintain a record of the test inputs for which $S$ is executed. We use this dataset to train a surrogate model $SM$ and compute its error $e$ (line 5):

2: $DS^l \leftarrow DS$; $flag \leftarrow \top$;
5: ($SM, e$) $\leftarrow$ Train($DS^l$);

Specifically, we train $SM$ using 80% of $DS^l$ and compute the mean absolute error of $SM$ on the remaining 20% of $DS^l$. The split ratio is based on the well-known 80/20 rule [80]. We employ a boolean variable, $flag$, to decide whether training a surrogate model is necessary. Initially, on line 2, we initialize $flag$ to $\top$ to ensure that a surrogate model is trained during the first iteration. Then, on line 7, we randomly generate a new test input and denote it by $t$:

7: $t \leftarrow$ GenerateTests($\mathcal{R}$, 1);

Then, on line 8, the surrogate model $SM$ predicts a fitness value for $t$. The predicted fitness value is denoted by $\bar{F}(t)$:

8: $\bar{F}(t) \leftarrow SM(t);$

On line 9, we calculate a confidence interval for $\bar{F}(t)$ based on the prediction error $e$. Specifically, we compute the interval $[\bar{F}(t) - e, \bar{F}(t) + e]$ as the confidence interval for $\bar{F}(t)$. If the confidence interval remains entirely within the positive or negative range, indicating a definite pass or fail verdict for the test input $t$, we skip executing system $S$ for the test input $t$. This is because we can confidently label it as either pass or fail. Subsequently, we add the test input $t$ along with its predicted fitness value (i.e., $\bar{F}(t)$) to $DS$ (line 10). In addition, we set $flag$ to $\bot$ in order to prevent retraining the surrogate model $SM$ in the next iteration (lines 11):

9: **if** $\exists \sim \in \{\geq, <\} \cdot (\bar{F}(t) \sim 0) \wedge (\bar{F}(t) \pm e \sim 0)$
10: $\quad DS \leftarrow DS \cup \{\langle t, \bar{F}(t) \rangle\};$
11: $flag \leftarrow \bot;$

Otherwise, on line 12, if the confidence interval spans both positive and negative numbers, indicating that it covers test inputs with both pass and fail verdicts, we proceed to execute system $S$ to calculate the actual fitness value for $t$. Then, we add $t$ along with its actual fitness value to $DS$ and $DS^l$ (lines 14-15). In this case, we set $flag$ to $\top$ to indicate that a system execution has taken place (line 15):

12: **else**
13: $\quad \{\langle t, F(t) \rangle\} \leftarrow \text{Execute}(S, \{t\}, F);$
14: $\quad DS^l \leftarrow DS^l \cup \{\langle t, F(t) \rangle\};$
15: $\quad DS \leftarrow DS \cup \{\langle t, F(t) \rangle\}; flag \leftarrow \top;$

In the latter case (i.e., lines 12-15), we re-train the surrogate model $SM$ using the dataset $DS^l$ which includes the new test input and its actual fitness value (lines 4-5). The algorithm returns the final dataset $DS$ when the execution budget runs out (line 18):

18: **return** $DS;$

The execution budget expires when either the system has been executed to the point where the size of $DS^l$ reaches its desired limit, or our time budget is exhausted, depending on which occurs first.

    In our experimentation (Section 4), we consider the surrogate-model types shown in Table 1. These surrogate-model types are the most widely used ones in the evolutionary search and software testing literature [43, 44, 78, 86]. As suggested by the literature and also as we show in our evaluation (Section 4.2), no surrogate-model type consistently outperforms the others [46, 92]. Therefore, it is recommended to use a combination of surrogate models [58]. In this paper, we propose, to our knowledge, a novel variation of Algorithm 2 where we train multiple surrogate models and use for predicting fitness values the model that has the lowest error. This variation is shown in Algorithm 3 where we change line 5 of Algorithm 2 to train and tune a list of surrogate models instead of just one model. We then select the surrogate model with the lowest error for making predictions until the next time we re-train the models. Similarly, each time we execute line 5 of Algorithm 2, we re-train a list of surrogate models and select the one that has the lowest error. We refer to our proposed variation as *dynamic* surrogate-assisted test generation.

    In both Algorithm 2 and the variation suggested in Algorithm 3, the first time we train a surrogate model, we also tune its hyperparameters using Bayesian optimization [84]. We use the same tuned hyperparameters in all future iterations. The cost of training and tuning surrogate models for the first time is on the same scale as the cost of a single execution of our CI systems. The time for

Table 1. Surrogate models and their descriptions.

| Name | Description | Name | Description |
|---|---|---|---|
| GL | Gaussian Process Regression – nonparametric Bayesian with linear kernel. | RT | regression tree. |
| GNL | Gaussian Process Regression – nonparametric Bayesian with nonlinear kernel. | RF | random forest. |
| LSB | Gradient Boosting – an ensemble of regression trees. | SVR | Support Vector Regression. |
| NN | a two-layer feedforward Neural Network. | | |

---

**Algorithm 3** Dynamically selecting surrogates

---

4: . . .
5: **for** i = 1 to $sm$ **do** //Train surrogate models $SM_1, ..., SM_{sm}$
6:   $(SM_i, e_i) \leftarrow \text{Train}(DS^l)$;
7: **end**
8: $(SM, e) \leftarrow$ Select $SM \in \{SM_1, \ldots, SM_{sm}\}$ with the lowest error $e$
9: . . .

---

subsequent re-training of surrogate models is nonetheless negligible since re-training does not involve any tuning. As we discuss in Section 4, the overhead of re-training surrogate models does not deteriorate performance compared to other alternatives.

*3.2.2 ML-Guided Test Generation.* ML-guided test generation uses ML models for identifying the boundary regions that discriminate pass and fail test inputs and iteratively concentrating test-input sampling to those regions. The idea is that, irrespective of the separability of the set of test inputs, ML models can shift the focus of sampling from the homogeneous regions where either fail or pass verdicts are scarce to regions where neither fail nor pass would be dominant. We consider two alternative ML models that can help us sample from such boundary regions: regression trees (Algorithm 4) and logistic regression (Algorithm 5). As we describe below, a regression tree approximates pass-fail boundaries in terms of predicates over inputs variables, while logistic regression infers a linear formula over input variables.

**Regression-Tree Guided Test Generation.** Algorithm 4 uses the *DS* dataset obtained from Algorithm 1 (the preprocessing phase) to train a regression-tree model (lines 1-3). In our regression-tree models, tree edges are labelled with predicates $v_i \sim c$ such that $v_i$ is an input variable, $c \in \mathbb{R}$ is a constant and $\sim \in \{\leq, >\}$. The tree leaves partition the given dataset into subsets such that information gain is maximized [91]. Each leaf is labelled with the average of the fitness measures of the test inputs in that leaf. Provided with a regression tree, Algorithm 4 identifies predicates $\{v_{i_1} \sim c_{i_1}, \ldots, v_{i_m} \sim c_{i_m}\}$ that appear on the two paths whose leaf-node values are closest to zero (one above and one below zero). These predicates specify the boundary between pass and fail, and, as such, we call them *boundary* predicates. By simplifying the boundary predicates, each variable can have at most one upper-bound predicate ($v \leq c$) and at most one lower-bound predicate ($v > c$). For each predicate $v_{i_j} \sim c$ where $\sim \in \{\leq, >\}$, the algorithm replaces the existing range $R_{i_j}$ of $v_{i_j}$ with $R'_{i_j} = [c-5\% \cdot c, c+5\% \cdot c]$ (lines 5-7). This will ensure that we sample $v_{i_j}$ within the 5% margin around the constant $c$. The variables that do not appear in the boundary predicates retain their range from the previous iteration. We note that if $R'_{i_j}$ does not reduce the range for $v_{i_j}$, i.e., the size of $R'_{i_j}$ is greater than $R_{i_j}$, we do not replace $R_{i_j}$ with $R'_{i_j}$. This is to ensure that larger ranges are not carried over to the next iteration if the range has already been narrowed at some previous iteration. Next, the algorithm generates a test input within the constrained search space (line 8), executes the test input, and adds it along with its fitness measure to *DS* (line 9). The algorithm returns the final dataset after the execution budget runs out (line 11). The execution budget expires when either the
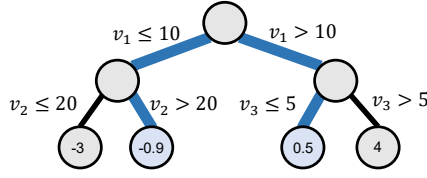
Fig. 4. A regression tree trained in an iteration of Algorithm 4. The figure illustrates two regression tree paths chosen for test generation on line 4 of Algorithm 4.

---

**Algorithm 4** ML-guided test generation with regression trees

---

**Input** $S$: System
**Input** $\mathcal{R} = \{R_1, \ldots, R_n\}$: Ranges for input variables $v_1$ to $v_n$
**Input** $F$: Fitness Function
**Input** INITIALDATASETSIZE: size of initial dataset
**Output** $DS$: A dataset to train failure models

1: $DS \leftarrow$ Preprocessing($S$, $\mathcal{R}$, $F$, INITIALDATASETSIZE);//Run Algorithm 1 to generate an initial dataset
2: **while** (execution budget remains) **do**
3:   $RegTree \leftarrow$ Train($DS$);
4:   Let $R'_{i_1}, \ldots, R'_{i_m}$ be reduced ranges obtained from $RegTree$;
5:   **for** each variable $v_{i_j}$ s.t. $j \in \{1, \ldots, m\}$ **do**
6:     $\mathcal{R} \leftarrow (\mathcal{R} \setminus \{R_{i_j}\}) \cup \{R'_{i_j}\}$; // Replace the range $R_{i_j}$ of $v_{i_j}$ in $\mathcal{R}$ with the new reduced range $R'_{i_j}$ from line 4
7:   **end**
8:   $\{t\} \leftarrow$ GenerateTests($\mathcal{R}$, 1); //Generate one test input
9:   $DS \leftarrow DS \cup$ Execute($S$, $\{t\}$, $F$); // Compute fitness for $t$ and add to $DS$
10: **end**
11: **return** $DS$

---

system has been executed to the point where the size of $DS$ reaches its desired limit, or our time budget is exhausted, depending on which occurs first.

To illustrate range reduction for variables using regression trees, consider the example in Figure 4. We choose the two thicker paths highlighted in blue since their leaf-node values are closest to zero. Variables $v_1$, $v_2$ and $v_3$ appear on these paths. Suppose the initial ranges for $v_1$, $v_2$ and $v_3$ to be $[0, 20]$, $[10, 30]$, and $[1, 7]$, respectively. Using the regression tree and the process described above, the new reduced ranges for $v_1$, $v_2$ and $v_3$ are $[9.5, 10.5]$, $[19, 21]$ and $[4.75, 5.25]$ respectively. By sampling within these ranges, we get to focus test-input generation on the pass-fail border identified by the regression tree.

**Logistic Regression Guided Test Generation.** Similar to Algorithm 4, Algorithm 5 uses the dataset $DS$ from the preprocessing phase to train a logistic regression model (lines 1-3). Since logistic regression is a classification technique, the quantitative labels in $DS$ are replaced with pass/fail labels before training. A linear logistic regression model is represented as $\log(\frac{p}{1-p}) = c + \sum_{i=1}^{n} c_i v_i$ where $v_1, \ldots, v_n$ are the input variables, $c_i$'s and $c$ are co-efficients, and $p$ is the probability of the pass class [11, 91]. The algorithm then randomly samples a few test inputs in the search space and picks the one closest to the logistic regression formula obtained by setting $p$ to the percentage of the pass labels in $DS$ (line 4-5).
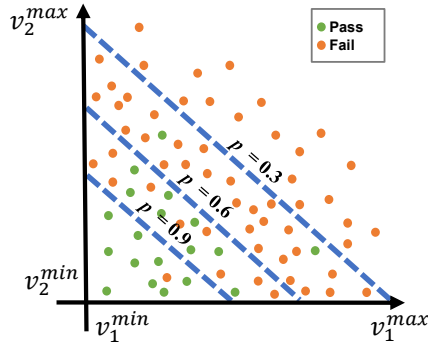
Fig. 5. Illustrating logistic regression lines for different values of $p$ assuming that we have two variables $v_1$ and $v_2$. The value of $p$ is used on line 5 of Algorithm 5 to generate a test close to the regression border.

---

**Algorithm 5** ML-guided test generation with logistic regression

---

**Input** $S$: System
**Input** $\mathcal{R} = \{R_1, \ldots, R_n\}$: Ranges for input variables $v_1$ to $v_n$
**Input** $F$: Fitness Function
**Input** INITIALDATASETSIZE: size of initial dataset
**Output** $DS$: A dataset to train failure models

1:   $DS \leftarrow$ Preprocessing($S$, $\mathcal{R}$, $F$, INITIALDATASETSIZE);//Run Algorithm 1 to generate an initial dataset
2:   **while** (execution budget remains) **do**
3:     $LogReg \leftarrow$ Train($DS$);
4:     $p \leftarrow$ Probability($DS$); // probability of pass in $DS$
5:     $t \leftarrow$ GenerateCloseToRegBorder($\mathcal{R}$, $LogReg$, $p$); // Select a test close to the regression border for $p$
6:     $DS \leftarrow DS \cup$ Execute($S$, $\{t\}$, $F$);// Compute fitness for $t$ and add to $DS$
7:   **end**
8:   **return** $DS$

---

By assigning a value to $p$ in the logistic formula above, the formula turns into a linear equation. Figure 5 shows examples of such linear equations for different values of $p$ assuming that we have two variables $v_1$ and $v_2$. The line with $p = 0.9$ identifies a region where the majority of test inputs pass. On the other end of the spectrum, the line with $p = 0.3$ identifies a region where the majority of test inputs fail. Setting $p$ to the percentage of pass in $DS$ is a heuristic to identify a region that includes a mix of pass and fail test inputs. Our sampling should, therefore, exploit this region. The test input selected on line 5 along with its fitness measure computed using $S$ is added to $DS$ (line 6). The algorithm re-trains the logistic regression model whenever a test input has been added to $DS$ (line 3). The algorithm returns the final dataset when the execution budget runs out (line 8). The execution budget expires when either the system has been executed to the point where the size of $DS$ reaches its desired limit, or our time budget is exhausted, depending on which occurs first.

Similar to Algorithm 2, the hyperparameters of the regression-tree and logistic-regression models are tuned using Bayesian optimization the first time the models are built, and the same tuned hyperparameters are used in all future iterations.

## 3.3 Building Failure Models

The output of the main loop in Figure 1 is a set *DS* of tuples $\langle t, F(t) \rangle$ where $t$ is a test input and $F(t)$ is its fitness value. We first convert *DS* into a dataset where test inputs are labelled by *pass* and *fail* labels. Provided with a labelled dataset, we use decision-rule models built using RIPPER [40] to train failure models. Decision-rule models generate a set of IF-condition-THEN-prediction rules where the condition is a conjunction of predicates over the input features and the prediction is either pass or fail. In Section 2, we already showed two examples of such rules for spurious failures. When no domain knowledge is available, one can directly use the input variables of the system ($S$) as features for learning. When domain knowledge is available, feature design for decision-rule models can be improved in two ways: (1) Excluding input variables that are orthogonal to the requirement under analysis. For example, the prerequisite for the requirement $\varphi$ in Section 2 is that the autopilot should be enabled, i.e., APEng = on. As far as test generation for $\varphi$ is concerned, we need to set APEng = on, since otherwise, $\varphi$ holds vacuously. We thus do not use APEng as an input feature. For another example, in $\varphi$, we do not use the desired altitude as an input feature either, since the system is expected to satisfy $\varphi$ for any desired altitude in the default range. (2) Using domain knowledge to formulate features over multiple input variables. For NTSS, as discussed in Section 2, the goal is to identify limits on the traffic that can flow through NTSS classes without compromising network quality. Based on domain knowledge, we know that flows have a cumulative nature. Hence, for NTSS, we use as features *sums of subsets* of flow variables. Naturally, like in any feature engineering problem, one can hypothesize alternative ways of formulating the features and empirically determine the formulation leading to highest accuracy [79].

## 4 EVALUATION

In this section, we evaluate our approach by answering the following research questions (RQs):

**RQ1 (Configuration).** *Which surrogate-assisted technique offers the best trade-off between accuracy and efficiency?* We compare *eight* surrogate-assisted algorithms. These eight algorithms are: (a) Algorithm 2 used with the seven surrogate models in Table 1 individually, and (b) the dynamic surrogate algorithm (Algorithm 3) that uses the seven surrogate models simultaneously and selects the best model dynamically. To measure accuracy, we check the correctness of the labels of the tests in the generated datasets; and, to measure efficiency, we evaluate the size of the generated datasets. We use the optimal algorithm for answering RQ2 to RQ4.

**RQ2 (Effectiveness).** *How accurate are the failure models generated by the surrogate-assisted and ML-guided techniques?* We evaluate and compare the accuracy of the failure models obtained by surrogate-assisted and ML-guided algorithms as well as those obtained based on randomly generated test inputs (random baseline).

**RQ3 (SoTA Comparison).** *How accurate are the failure models generated in RQ2 compared to those generated by the state of the art (SoTA)?* We use the top-performing technique from RQ2 to compare against SoTA. Among the existing approaches that build failure-inducing models [51, 65, 66], we select the Alhazen framework [65], since it uses interpretable machine learning. While Alhazen is geared towards systems with structured inputs (as opposed to systems with numeric inputs, i.e., the focus of our work), in the absence of baselines for systems with numeric inputs, Alhazen is our best baseline for comparison. To be able to compare with Alhazen, we adapt it to numeric-input systems, as we describe in Section 4.4.

**RQ4 (Usefulness).** *How useful are failure models for identifying spurious failures?* We answer this question for the most accurate failure models from RQ2 and for the two CI systems, namely NTSS and autopilot, discussed in Section 2. NTSS and autopilot are representative examples of industrial systems in the network and CPS domains, respectively. For both systems, we validate the

Table 2. Names, descriptions, the number of requirements and the identifiers of our study subjects. For each subject, we indicate if it is computer-intensive (CI). All artifacts including requirements statements are available in our supplementary material [6].

| Name | Description | #Reqs | ID | CI |
|---|---|---|---|---|
| Tustin | A common flight control utility for computing the Tustin Integration – A Simulink model with 57 blocks. | 9 | TU1...TU9 | ✗ |
| Regulator | A regulators inner loop architecture used in many feedback control applications – A Simulink model with 308 blocks. | 1 | REG | ✗ |
| Nonlinear Guidance | A nonlinear algorithm for generating a guidance command for an air vehicle – A Simulink model with 373 blocks. | 1 | NLG | ✗ |
| Finite State Machine | A finite state machine to enable autopilot mode if a hazardous situation is identified – A Simulink model with 303 blocks. | 1 | FSM | ✗ |
| Autopilot | A single-engine, high-wing, propeller-driven aircraft simulation with all six degrees of freedom – A Simulink model with 1549 blocks. | 3 | AP1, AP2, AP3 | ✓ |
| Network Traffic Shaping System | An NTSS testbed [64] developed using three virtual machines and based on OpenWRT [12]. | 1 | NTSS | ✓ |

failure-inducing rules against domain knowledge to determine whether the resulting failures are genuinely spurious.

## 4.1 Study Subjects

Our study subjects, which are listed in Table 2, originate from the network and CPS domains. Below, we introduce our study subjects and discuss how these subjects satisfy assumptions **A1** and **A2** provided at beginning of Section 3.

*Network-system subject.* Our network-system subject is the Network Traffic Shaping System (NTSS) discussed in Section 2. To test NTSS, we transmit flows with different bandwidth values into different NTSS classes. A test input for NTSS is defined as a tuple $t = (v_1, \ldots, v_n)$ where $n$ is the number of NTSS classes, and each variable $v_i$ represents the bandwidth of the data flow going through class $i$. The fitness function for NTSS measures the network quality based on the well-known mean opinion score (MOS) metric [85]. This fitness function ensures assumption **A2** [64]. For our experiments, we use an NTSS setup based on an industrial small-office and home-office use case from our earlier work [64]. This setup runs Common Applications Kept Enhanced (CAKE) [57], which is an advanced and widely used traffic-shaping algorithm. The setup uses the 8-tier mode of CAKE known as diffserv8 [27, 57], i.e., the number of NTSS classes is 8.

*Simulink subjects.* Simulink [38] is a widely used language for specification and simulation of CPS. The inputs and outputs of a Simulink model are represented using signals. A typical input-signal generator for Simulink characterizes each input signal using a triple $(int, R, n)$ such that $int$ is an interpolation function, $R$ is a value range, and $n$ is a number of control points [29, 87]. Let $(x, y)$ be a control point. The value of $x$ is from the signal's time domain, and the value of $y$ is from range $R$. Provided with $n$ control points, the interpolation function $R$ (e.g., piecewise constant, linear or piecewise cubic) constructs a signal by connecting the control points [87]. It is usually assumed that the input signals for a Simulink model have the same time domain. Further, for the purpose of testing, we make the common assumption that the control points are equally spaced over the time domain, i.e., the control points are positioned at a fixed time distance [29, 47, 87]. Hence, to generate test inputs, we only need to vary the $y$ variable of control points in the range $R$. Note that the type of the interpolation function is fixed for each Simulink-model input as the interpolation

function type is determined by the meaning of the input signal. For example, a reference signal is often a constant or step function. As a result, we can exclude the interpolation function from the test-input signals, and define each test-input signal as a vector of $n$ control variables taking their values from $R$. Simulink models often have clearly stated requirements. To define a fitness function for each requirement, we encode the requirement in RFOL – a variant of the signal temporal logic which is expressive enough for capturing many CPS requirements [74]. For our fitness function, we utilize the RFOL semantic function – a quantitative measure that conforms to assumption **A2** as shown by Menghi et al. [74].

In our evaluation, we use a public-domain benchmark of Simulink specifications from Lockheed Martin [10]. This benchmark provides a set of representative CPS systems that are shared by Lockheed as verification-and-validation challenge subjects for researchers and quality-assurance tool vendors. The benchmark is comprised of eleven specifications, of which six were not useful for our evaluation. These six specifications had requirements that either did not fail, or failed for all test inputs and thus, their failure models could be trivially defined as the entire input space. In our experiments, we focus on five of the Simulink specifications from Lockheed's benchmark. The first five rows of Table 2 list these specifications that have a total of 15 requirements combined.

In total, we have 16 requirements: one for NTSS, and 15 for the five Simulink specifications listed in Table 2. As discussed in Section 3, we define one fitness function per requirement. Hence, in total, we have 16 different subjects for our evaluation. Among these, four are compute-intensive (CI) and twelve are non-CI. Both NTSS and autopilot are CI: On average, each execution of NTSS takes $\approx 4.5min$, and each execution of autopilot takes $\approx 0.5min$. The execution times for non-CI subjects are negligible ($< 1s$). All experiments were conducted on a machine with a 2.5 GHz Intel Core i9 CPU and 64 GB of DDR4 memory.

## 4.2 RQ1-Configuration

We compare eight versions of the surrogate-assisted algorithm. Seven are Algorithm 2 used with an individual surrogate model from those in Table 1. We refer to each of these algorithms as SA-XX where XX is the name of the surrogate model from Table 1. For example, SA-NN refers to Algorithm 2 used with NN. The final (i.e., eighth) algorithm is the dynamic surrogate-assisted one (Algorithm 3). We refer to Algorithm 3 as SA-DYN.

For RQ1, we use the 12 non-CI subjects in Table 2. Performing RQ1 experiments on CI subjects would be prohibitively expensive. For example, an approximate estimate for the execution time of the experiments required to answer RQ1 is over a year, if the experiments are performed on NTSS using the same experimentation platform. Thus, we opt for the non-CI subjects for RQ1.

*Setting.* For each subject, we run the eight SA-XX algorithms for an equal time budget. The time budget given for each subject depends on the subject execution time. The detailed time budgets are available in our supplementary material [22]. To account for randomness, we repeat each algorithm for each subject ten times.

*Metrics.* Recall that the output of the main loop in Figure 1 for surrogate-assisted algorithms is a set *DS* where the test inputs are labelled with either predicted or actual fitness values. To measure efficiency, we take the cardinality of the generated dataset, *DS*. Since the algorithms have the same time budget, an algorithm is more efficient than another if it generates a larger dataset. To construct failure models, we classify test inputs as pass or fail based on their fitness values. A dataset is accurate if it contains few test inputs with *incorrect labels*, i.e., test inputs with inconsistent pass/fail labels based on their predicted versus actual fitness values. To obtain the actual fitness values for all test inputs, we run the system using those test inputs for which only surrogate-assisted algorithms provided predicted fitness values. Note that this step is intended exclusively for our experiments and is not a component of our main approach. Next, we count the number of tests in which the
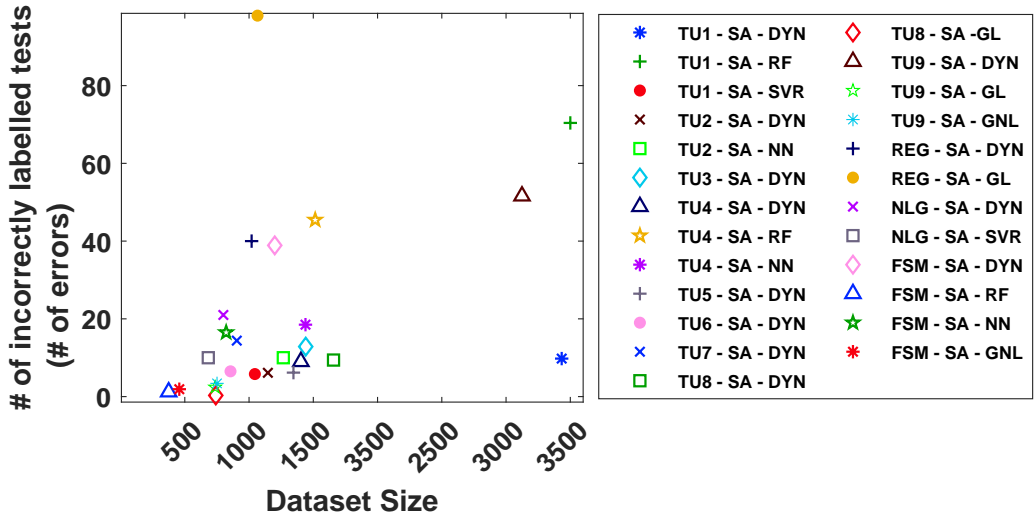
Fig. 6. Comparing datasets generated by eight different surrogate-assisted algorithms with respect to the number of errors in the datasets and the dataset size.
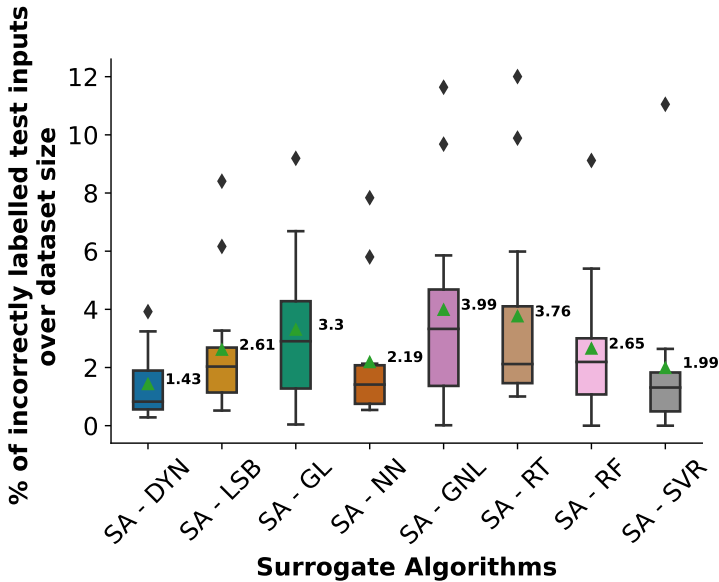


Fig. 7. Percentages of the incorrectly labelled tests over the dataset size for different surrogate-assisted algorithms.

predicted label differs from the actual one. The number of incorrectly labelled tests serves as a measure of error or inaccuracy for the surrogate-assisted algorithms.

*Results.* The scatter plot in Figure 6 shows the results of the experiments for RQ1. The x-axis indicates $|DS|$ and the y-axis indicates the number of incorrectly labelled tests in $DS$. Each point shows the result of applying one SA-XX algorithm to one subject. The 12 subjects are indicated

by TU1 to TU9, REG, NLG, and FSM (see Table 2). For each subject, an algorithm is considered better when it generates larger datasets with fewer errors. Since we compare eight algorithms for 12 subjects, we would need 96 points to show all the results. To reduce clutter, for each subject, we only show the Pareto Front (PF) points. That is, for each subject, we only show the algorithms that are not dominated by other algorithms either in terms of the number of errors or in terms of the dataset size. For example, for subject TU1, algorithms SA-DYN, SA-RF and SA-SVR dominate other algorithms and offer the best trade-off between the number of errors and the dataset size. As Figure 6 shows, for four subjects, SA-DYN is the only best trade-off (i.e., PF point), and for eight subjects, it is one of the best PF points. For the latter eight cases, SA-DYN offers an alternative that, compared to the other PF points, either has considerably fewer errors while its dataset is not much smaller, or its dataset is considerably larger while the number of errors is not much higher. For example, for TU1, SA-DYN, compared to SA-RF, provides 60 less incorrectly labelled tests, while the dataset sizes are almost the same (3433 for SA-DYN vs. 3500 for SA-RF). Also, compared to SA-SVR, SA-DYN provides a larger dataset (3433 vs. 1045) while the number of incorrectly labelled tests is almost the same (10 vs. 6).

Figure 7 shows the ratios of the number of errors (i.e., the number of incorrectly labelled tests) over $|DS|$ for different SA algorithms and for all the 12 subjects. The SA-DYN algorithm has the lowest average error which is 28% lower than that of the second best algorithm, SA-SVR, i.e., $\frac{1.99-1.43}{1.99} = 28\%$. We compare the results in Figure 7 using the non-parametric pairwise Wilcoxon rank-sum test [72] and the Vargha-Delaney's $\hat{A}_{12}$ effect size [88]. The SA-DYN algorithm is statistically significantly better than other algorithms with a high effect size for GL, GNL and LSB, a medium effect size for RT, a small effect size for NN and RF, and a negligible effect size for SVR. The comparison of the dataset sizes shows that SA-DYN generates datasets that are significantly larger than those generated by SVR with a large effect size. Further, SA-DYN generates significantly larger datasets that are at least 33% larger than those obtained from other algorithms. Figures comparing dataset sizes and statistical tests for RQ1 are available in our supplementary material [9, 23].

> Our evaluation for RQ1 performed using seven surrogate-model types in the literature (see Table 1) shows that, compared to using surrogate models individually, our dynamic surrogate-assisted approach provides the best trade-off between accuracy and efficiency.

## 4.3 RQ2-Effectiveness

We compare SA-DYN, i.e., the best approach identified in RQ1, with the two ML-guided algorithms described in Section 3.2.2 as well as a standard adaptive random-search algorithm. In the remainder of this section, we refer to SA-DYN as SA. We use RT and LR to refer to the ML-guided algorithms that employ regression tree (Algorithm 4) and logistic regression (Algorithm 5), respectively. We use RS for adaptive random search.

*Setting.* We apply the four algorithms (i.e., SA, RT, LR, and RS) to the 16 subjects in Table 2. For each CI subject, we execute the four algorithms for an equal time budget and then compare the results. For the non-CI subjects, however, fixing the time budget may favour RS, as the other three algorithms have an additional overhead for training ML models. This overhead time is negligible when compared to the execution time for CI subjects. But, for non-CI subjects, we can execute many tests within the overhead time, which can skew the results in favour of RS. Therefore, to ensure that the results are valid for CI subjects, we follow the approach proposed by Menghi et al. [73]. Specifically, we use the execution time of CI subjects to limit the number of test inputs that each algorithm executes for non-CI objects. Briefly, to compare two algorithms with different overhead times, we allow the algorithm with the lower overhead time to execute $x$ additional test

inputs such that $x$ multiplied by the execution time of a typical CI subject (instead of a non-CI subject) is equal to the difference in the two algorithms' overhead time. For the non-CI Simulink subjects in Table 2, we use the average execution time of the Simulink CI subject, i.e., autopilot. Given a time budget, we compute the maximum number of test executions that each of the SA, LR, RT, and RS algorithms can perform within this time budget for autopilot. We then use these numbers to cap the number of test executions for each algorithm when we compare them for the non-CI models in Table 2. The time budget we consider for comparing these algorithms for CI subjects and the maximum number of test executions we use to compare them for non-CI subjects are available in our supplementary material [24, 25]. We repeat each algorithm twenty times for each subject except NTSS. For NTSS, due to its large execution time, we repeat each algorithm only ten times.

*Metrics.* We use the datasets created by SA, LR, RT, and RS to build decision rule models (DRM). For hyperparameter tuning, we use Bayesian optimization [84]. To avoid bias towards any particular algorithm, we use for tuning the union of the datasets obtained from our four algorithms. Having fixed the hyperparameters, we train a DRM separately for each dataset obtained from each repetition of our four algorithms. We evaluate the DRMs using three metrics: *accuracy*, *precision* and *recall*. Accuracy is the number of correctly predicted tests over the total number of tests. Since DRMs are mainly used to predict the failure class, we compute precision and recall for the failure class as follows: Precision is the number of fail-class predictions that actually belong to the fail class, and recall is the number of fail-class predictions out of all the actual failed tests in the dataset. We use randomly generated test inputs within the variables' default ranges to measure the accuracy, precision and recall of each DRM.

*Features for learning.* For the Simulink subjects, we use as features the individual input variables of each subject model but exclude the following two kinds of variables: (1) Variables that are explicitly fixed to a value in a requirement (e.g., variable APEng discussed in Section 3.3). (2) Reference variables that indicate the desired value of a controlled process, noting that the system is expected to satisfy its requirements for *any* value in a reference variable's valid range. As such, reference variables cannot contribute to creating failure conditions. The desired altitude variable discussed in Section 3.3 is an example of a reference variable.

For NTSS, as discussed in Section 3.3, we consider alternative features as follows: the set of all individual variables (i.e., individual NTSS classes), sums of two variables, sums of three variables, ..., and the sum of all eight variables. We then create, for each input feature, one DRM based on a dataset obtained by each of the four algorithms. In total, for each algorithm, we create 248 DRMs for NTSS. That is, the sets of all subsets larger than two (247) and the set of all individual variables. Given the large number of hypothesized input features, we evaluate the accuracy of the resulting DRMs and keep the input features that yield reasonably high accuracy over multiple runs of SA, LR, RT, RS. This results in the retention of two input features for NTSS with an accuracy higher than 80%.

*Results.* Figures 8(a)-(c) compare across all the 16 subjects the average accuracy of DRMs obtained by SA, LR and RT (on y-axis) against the average accuracy of DRMs obtained by RS (on x-axis). Each point in each of Figures 8(a)-(c) corresponds to one study subject. A blue point indicates that the difference in accuracy is statistically significant as per the Wilcoxon rank-sum test. The DRMs obtained using SA are significantly more accurate than those obtained using RS for 14 out of 16 subjects, including all the CI subjects. The accuracy of the DRMs obtained using LR is significantly better than that obtained using RS for seven subjects. The accuracy of the DRMs obtained using RT is significantly better than that obtained using RS for nine subjects. Overall for all the subjects, SA has the highest average accuracy (83%), followed by RT and LR with average accuracies of 78% and 76%, respectively. RS has the lowest average accuracy (72%). Finally, the accuracy of SA, RT and LR
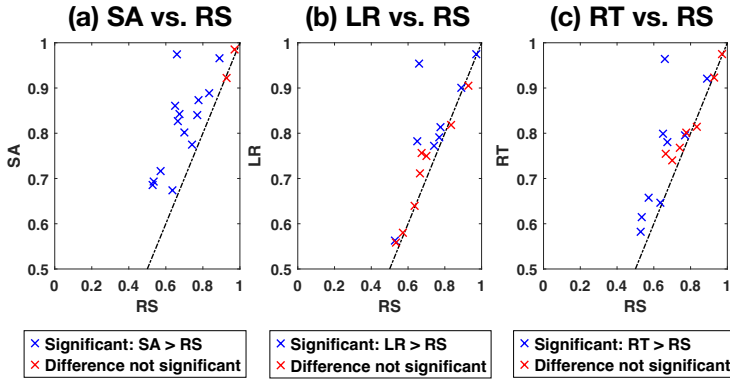
Fig. 8. Comparing the accuracy of decision-rule models obtained based on the datasets generated by SA, LR, and RT against those obtained by RS for all the 16 subjects in Table 2.
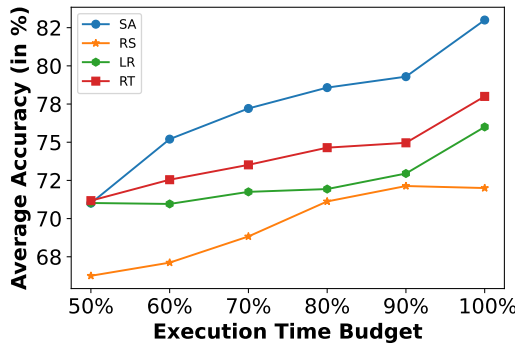


Fig. 9. Average accuracy of the DRMs obtained using SA, RL, RT and RS for the 16 subjects as the time budget is varied.

is significantly better than that of RS. The effect size for SA versus RS is medium, and the effect size for RT and LR versus RS is small.

Figure 9 compares the average accuracy of the DRMs obtained using SA, LR, RT and RS for our 16 subjects as the execution-time budget is varied from 50% to 100%. Note that in our experiments, we dedicate 50% of the time budget to the prepossessing step. Therefore, Figure 9 compares the impact of the four algorithms over the remaining 50% of the execution-time budget. As the figure shows, SA consistently has the highest average accuracy. Further, the average accuracy of RS reaches a plateau, whereas the other three algorithms keep improving as the budget increases. The main reason for this difference is that RS, unlike the other algorithms, does not use machine learning models to guide its test input sampling. Specifically, RS generates test inputs randomly across the entire input space. In contrast, RT and LR exploit boundary regions that separate passing and failing test inputs, resulting in a steady increase in accuracy as the search budget grows. Further, SA utilizes surrogate models, generating significantly more test inputs compared to the other algorithms within the same allotted time budget. As a result, SA achieves the highest average accuracy among all the algorithms.
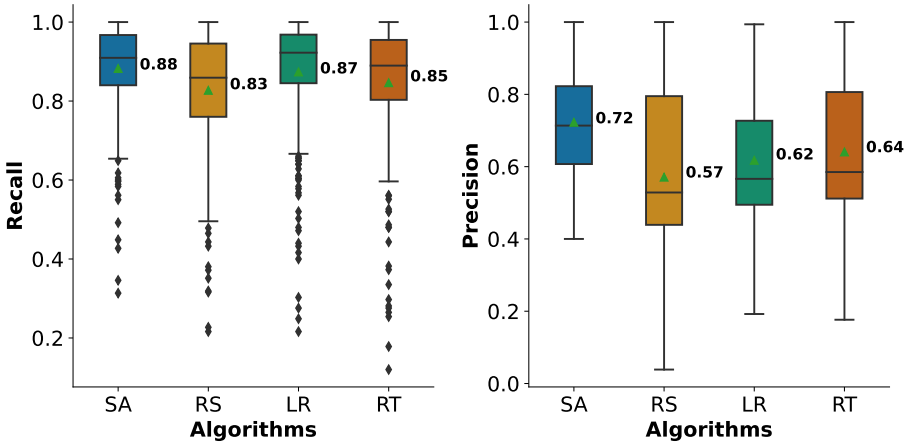
Fig. 10. Recall and Precision for the DRMs obtained based on SA, RL, RT and RS for all the 16 subjects in Table 2.

Figure 10 compares the recall and precision for all the DRMs obtained using SA, LR, RT and RS for our 16 subjects. Recall measures the ability of DRMs to precisely identify the failure conditions, whereas precision assesses the ability of DRMs to generate failure instances correctly. The SA algorithm has the highest average recall (88%), followed by LR (87%) and RT (85%). RS has the lowest average recall (83%). Moreover, SA has the highest average precision (72%), followed by RT (64%) and LR (62%) and RS has the lowest average precision (57%). The recall of SA, RT and LR is significantly better than that of RS with a medium effect size for SA, a small effect size for LR, and negligible effect size for RT. Likewise, the precision of SA, RT and LR is significantly better than that of RS with a medium effect size for SA, small effect size for RT and negligible effect size for LR. Finally, the accuracy, recall and precision values for SA are significantly higher than those for RT and LR. Precision and recall for SA exhibit similar trends to that for accuracy (Figure 9). Detailed charts comparing precision and recall for the four algorithms as the execution-time budget varies are available online [8, 19, 26].

> Our evaluation for RQ2 performed over all our study subjects shows that our dynamic surrogate-assisted approach yields failure models with significantly higher accuracy, precision and recall compared to those obtained using ML-guided algorithms and a random baseline.

## 4.4 RQ3-SOTA Comparison

We compare SA – the top-performing algorithm from RQ2 – with an adaptation of Alhazen to numeric-input systems. We hereafter use SoTA to refer to this adaptation, discussed next. Similar to SA, in SoTA, we generate test inputs according to the description in Section 4.1. The workflow of SoTA then matches the steps in Figure 1 with the difference that the model trained and refined in the main loop (i.e. step 3) is always a decision tree model; this decision tree is returned as the failure model at the end. Recall that our approach has a separate step, i.e., step 5, after the main loop to build failure models from the data generated by different algorithms. This step is not required in SoTA, noting that the model that SoTA refines during its main loop is used as the failure model. Algorithm 6 shows our implementation of SoTA. SoTA starts from an initial dataset (line 1). In each iteration, it builds a decision tree on the dataset (line 3). It then generates test inputs using all the

---

**Algorithm 6** Our implementation of SoTA

---

**Input** $S$: System
**Input** $\mathcal{R} = \{R_1, \ldots, R_n\}$: Ranges for Input variables $v_1$ to $v_n$
**Input** $F$: Fitness Function
**Input** INITIALDATASETSIZE: size of initial dataset
**Output** *DecisionTree*: A decision tree model (failure model)

1: $DS \leftarrow$ Preprocessing($S, \mathcal{R}, F$, INITIALDATASETSIZE); //Run Algorithm 1 to generate an initial dataset
2: **while** (execution budget remains) **do**
3:     *DecisionTree* $\leftarrow$ Train($DS$);
4:     Let $P_1, \ldots, P_q$ be all the paths obtained from *DecisionTree*; // Extract all the paths from the decision tree
5:     **for** each path $P_k$ s.t. $k \in \{1, \ldots, q\}$ **do** // Generate a test in each path based on the ranges obtained from that path
6:         Let $R'_{i_1}, \ldots, R'_{i_m}$ be reduced ranges obtained from $P_k$;
7:         **for** each variable $v_{i_j}$ s.t. $j \in \{1, \ldots, m\}$ **do**
8:             $\mathcal{R} \leftarrow (\mathcal{R} \setminus \{R_{i_j}\}) \cup \{R'_{i_j}\}$; // Replace the range $R_{i_j}$ of $v_{i_j}$ in $\mathcal{R}$ with the new reduced range $R'_{i_j}$ from line 6
9:         **end**
10:         $\{t\} \leftarrow$ GenerateTests($\mathcal{R}$, 1); // Generate a test in path $P_k$
11:         $DS \leftarrow DS \cup$ Execute($S, \{t\}, F$); // Compute fitness for $t$ and add to $DS$
12:     **end**
13: **end**
14: **return** *DecisionTree*

---

paths in the decision tree (lines 4-10). These test inputs are executed and added to the dataset along with their labels (line 11). The final decision tree is returned on line 14.

*Setting.* For this comparison, we apply SoTA to our CI-subjects in Table 2, i.e., NTSS, AP1, AP2 and AP3. For the decision tree parameters, e.g. maximum depth of tree and class weight, we use the same parameters as in the original study [21, 65]. We execute SoTA for the same time budget as SA. For SA, we use the dataset generated in RQ2. We repeat SoTA for twenty times for each subject except for NTSS. For NTSS, we repeat it only ten times due to the expensive execution time.

*Metrics.* In order to compare SoTA and SA, we build decision trees based on the datasets generated by SA in RQ2. To do so, we use the same decision tree parameters as those used by SoTA. We compare the decision trees using the three metrics explained in RQ2, i.e. accuracy, precision for fail class and recall for fail class. We also use the same test set utilized in RQ2.

*Results.* Figure 11 compares the average accuracy of the decision trees obtained by SA (on y-axis) against those obtained by SoTA (on x-axis) across the four CI subjects in Table 2. Similar to Figure 8, each point on Figure 11 corresponds to one study subject and a blue point denotes a statistically significant difference in accuracy, determined using the Wilcoxon rank-sum test. As figure 11 shows, the decision trees obtained using SA are significantly more accurate than those obtained by SoTA, for three out of the four subjects with a medium effect size. For the fourth subject there is no statistically significant difference.

Figure 12 compares the average accuracy of the decision trees obtained using SA and those obtained by SoTA for the four CI subjects as the execution-time budget is varied from 50% to 100%.
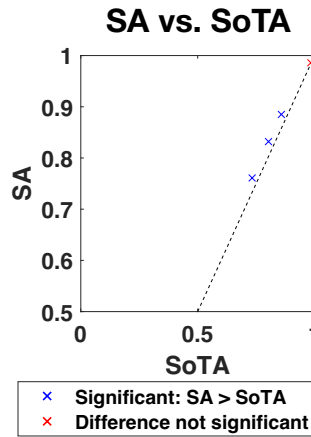
Fig. 11. Comparing the accuracies of the decision trees obtained on the datasets generated by SA and the decision trees returned by SoTA for all the four CI subjects in Table 2.

As the figure shows, the average accuracy of SA is consistently higher than SoTA as the time budget increases.

Finally, Figure 13 compares the recall and precision for all the decision trees obtained using SA and those obtained by SoTA for the four subjects. As shown by the figure, the average recall of SA (85%) is higher than SoTA (83%). Further, the recall of SA is significantly better than SoTA with small effect size. Moreover, the average precision of SA (76%) is higher than SoTA (73%). Similar to recall, the precision of SA is significantly better than SoTA with small effect size.

> Our evaluation for RQ3 performed using dynamic surrogate-assisted and the state-of-the-art approach over CI subjects indicates that our dynamic surrogate-assisted approach yields failure models with higher accuracy, precision and recall compared to those obtained from the state-of-the-art approach.

## 4.5 RQ4-Usefulness

In view of the results of RQ2, we use the DRMs generated by the SA algorithm to evaluate their usefulness in identifying spurious failures. We focus on the DRMs for the CI subjects in Table 2, i.e., NTSS, AP1, AP2 and AP3. For these subjects, we can validate whether the inferred rules lead to genuinely spurious failures. For NTSS, we had access to an expert from industry, and for autopilot, we had detailed requirements and design documents. Recall from Section 3.3 that the rules we obtain from DRMs are in the form of IF-condition-THEN-prediction. Each rule has a confidence that shows what percentage of the tests satisfying the condition of the rule conform to the rule's prediction label. From the DRMs for each of the four CI subjects, we extract the rules that predict the fail class with a confidence of 100%. These rules are candidates for specifying spurious failures, since they identify conditions that lead to and only to failures. We then select the rules that are not subsumed by others through logical implication. We use the Z3 SMT solver [41] to find logical implications. In the end, we obtain seven rules for NTSS, 17 rules for AP1, 24 rules for AP2 and 15 rules for AP3. On average, the rules for NTSS include two variables and three predicates, and the rules for autopilot include three variables and four predicates.

To validate the rules for NTSS, we presented the rules to a domain expert. Our domain expert for NTSS is a seasoned network technologist and software engineer with more than 25 years of
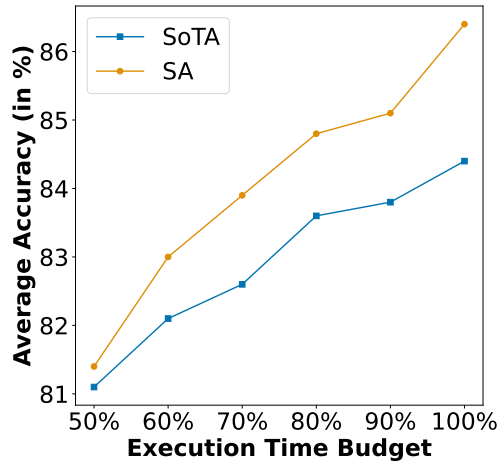
Fig. 12. Average accuracy of the decision trees obtained using SA and SoTA for the four CI subjects as the time budget is varied.
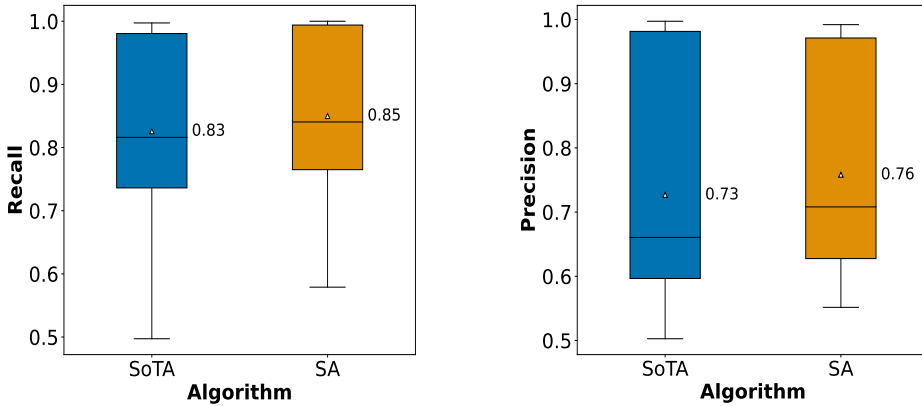


Fig. 13. Recall and Precision for the decision trees obtained based on SA and SoTA for four subjects in Table 2.

experience. The expert has been using the core enabling technology of NTSS (CAKE [57], discussed earlier in this section) in commercial networking solutions since 2018. Among the seven rules for NTSS, one rule constrains an input feature formulating a sum of the NTSS input variables. The rest of the rules involve predicates over individual input variables. The domain expert approved all the rules as they all correspond to situations where NTSS is overloaded with large traffic volumes, and hence, poor network quality is expected. We note that, for each input variable of NTSS, there is a threshold that is determined by the NTSS setup configuration. Among the seven rules, six constrain input variables near these known threshold values. An example of such rules is: class6 $\geq$ 91%· thresh6 $\wedge$ class7 $\geq$ 87%· thresh7 THEN FAIL where thresh6 and thresh7 are thresholds of class6 and class7, respectively. These rules matched the expert's intuition; nonetheless,

the expert still found the rules helpful as they provide data-driven evidence for what the expert could only estimate based on experience and ad-hoc observations rather than systematically collected data. More importantly, the rule constraining a sum of input variables (as discussed in Section 2) was of particular interest. For this rule, neither could the domain expert estimate the combination of the variables in the rule nor was the limit in the rule close to known thresholds. This rule prompted an investigation into the source code of CAKE. This investigation confirmed that classes 7, 6 and 5 have higher priority than other classes. As such, the combined cumulative usage of these three classes needs to be capped to maintain high quality of experience. The expert indicated that, without our approach, he would not have been able to formulate such a rule merely based on his existing test scenarios and expert knowledge.

For the autopilot's three requirements (AP1, AP2 and AP3) , we could conclusively confirm that 47 out of the 56 inferred rules represent spurious failures as per the handbook of the De Havilland Beaver aircraft [31]. For the remaining nine rules, we could not ascertain whether they represent spurious failures. These rules either fail due to real faults in the system or they are spurious, but further expertise is required to confirm their spuriousness. The full set of rules and our analysis are in our supplementary material [20].

Our results indicate that none of the 47 rules confirmed as inducing spurious failures for autopilot could be obtained solely from the datasets generated by the preprocessing phase of SA (see Figure 1). Similarly, only two of the seven rules for NTSS could be obtained using the datasets from the preprocessing phase. Furthermore, from the preprocessing phase datasets alone, no additional candidate rules could be inferred for spurious failures in either NTSS or autopilot. These findings highlight the importance of utilizing surrogate-assisted test generation in order to obtain useful and more comprehensive rules for spurious failures.

> Our validation of failure-inducing rules against domain knowledge indicates that our dynamic surrogate-assisted approach is effective for identifying spurious failures. Indeed, our results show that expert judgement alone or tests generated without the assistance of surrogates would miss many rules that one would be able to identify using our proposed approach.

## 4.6 Threats to Validity

The most important threats concerning the validity of our experiments are related to the internal and external validity.

*4.6.1 Internal Validity.* The *internal validity* risks are related to confounding factors. The effectiveness of failure-inducing rules inferred by our approach depends on the accuracy of fitness functions and the quality of the input datasets. For Simulink models, we use an automated and provably sound technique to obtain fitness functions from logical specifications [74]. However, the translation of natural-language requirements into logical specifications remains a manual task and necessitates domain-expert validation. In our experiments, we mitigated the risks associated with the accuracy of fitness functions as follows: For Simulink models, the fitness functions are automatically obtained from logical specifications approved by the engineers who developed the benchmark Simulink models [77]. For the NTSS case study, we validated the fitness function with our domain expert [64]. As for the risks related to the quality of datasets, we note that the labels for the data points are computed based on the actual system outputs, and hence, are always accurate. Further, we used adaptive random testing in the preprocessing step (see Figure 1) to diversify the generation of the datasets in the search input space.

To address any concerns regarding our comparison with SoTA in RQ3, we conducted a comprehensive review of the SoTA code [14] to ensure that the workflow of Algorithm 6, our adaptation of

SoTA, conforms to the original SoTA code. In addition, we employed the same hyper-parameters for decision trees in Algorithm 6 as those used by SoTA. Finally, we have disclosed the code of Algorithm 6 in our GitHub repository [4, 5] to facilitate further replication and comparison efforts.

*4.6.2 External Validity.* The subjects we used for our evaluation and the characteristics of these subjects may influence the generalizability of our results. Related to this threat, we note that: First, the Simuilnk models in Lockheed's benchmark represent realistic and representative CPS components from different domains. This benchmark has been previously used in the literature on testing CPS models [50, 77]. Second, our case studies are drawn from two different domains: CPS and networks. Third, our network case study (NTSS) represents an industrial system for which we could interact with a domain expert. The above being said, our work would benefit from further experiments with a broader class of systems.

## 5 LESSONS LEARNED

*Lesson 1. Decision rules are a better choice than decision trees for building failure models.* In this paper, we focused on interpretable ML techniques for building failure models. Using these techniques, we were able to generate constraints on system inputs that are easily understandable to humans [76]. Among interpretable ML techniques, decision rules and decision trees have been previously used in the literature for inferring rules pertaining to a specific behaviour of a system [37, 47, 48, 53, 65]. We chose to use decision-rule models in our work for the following reasons: Decision rules are known to produce fewer and more concise rules compared to decision trees, which often generate many rules involving several variables and predicates. Further, decision trees are prone to the replicated subtree problem [91]. This problem arises when the same subtree, with identical predicates and splits, appears multiple times in the tree. Replicated subtrees can increase model complexity, lead to overfitting, and hinder interpretability. Decision rules do not typically suffer from this problem, thus generally yielding more interpretable and less redundant rules. As mentioned in Section 4.5, the rules we obtain for NTSS and autopilot, on average, have three and four predicates over two and three variables, respectively. While one could argue that limiting the depth of a decision tree, as done by the SoTA baseline, would result in a reasonably small tree, our findings indicate that, decision trees built using the parameters of the SoTA baseline lead to a 40% higher number of rules and 10% more predicates compared to decision rules.

*Lesson 2. To evaluate a test generation algorithm for systems with numeric inputs, the accuracy and usefulness of the failure models produced by the algorithm offer more realistic insights about the algorithm than the number of individual failures found by the algorithm.* For systems with numeric inputs, slight modifications to the inputs of a failure-revealing test may lead to redundant failures, i.e., failures caused by the same fault. Even when one considers input diversity, e.g., measured by the Euclidean distance between test-input vectors, one cannot determine whether failures are non-redundant or valid by merely analyzing individual test inputs. Consequently, evaluating testing algorithms solely based on their ability to generate failures may result in misleading conclusions. Indeed, had we premised our evaluation on the number of detected failures, we would have inferred that our dynamic surrogate-assisted algorithm produces 2.3 times more failures compared to the ML-guided and the baseline algorithms. While this conclusion would strongly favour our approach, we do not believe that this large margin is an accurate representation of the degree of improvement that our algorithm delivers. Based on the results of RQ2 and RQ3, the dynamic surrogate-assisted algorithm, when compared to alternatives, leads to an accuracy improvement ranging from 2% to 14%.

## 6 RELATED WORK

Below, we discuss the related work on (a) the applications of ML in automated testing, (b) generating failure-inducing rules and (c) test input validation.

**Applications of Machine learning (ML) in automated testing.** ML has been widely used to enhance the effectiveness of fuzz testing [75] and search-based testing (SBT) [55]. In fuzz testing, ML has been employed to improve, among other things, seed generation, test sampling, and mutation-operator selection [54, 90]. In SBT, surrogate models developed based on ML have been used to effectively and efficiently test CI systems such as cyber-physical-system controllers and simulators [30, 60, 71, 73], and autonomous-driving systems [34, 35, 61]. These approaches demonstrate that using ML can improve the ability and the efficiency of testing in revealing faults. The ultimate goal of these approaches is to generate specific test cases. As such, these approaches are evaluated based on the number of failure-revealing tests and the severity of the failures, as determined by the fitness-function values.

Recent studies suggest that the focus of SBT should shift from generating a limited number of specific test cases to learning models that can explain system failures [45]. These models can then be employed for generating multiple test cases with specific properties. Motivated by these observations, our goal is to learn failure models and focus on improving the accuracy of these models for identifying spurious failures, rather than maximizing the number and severity of failure-revealing tests, which may not accurately reflect the context where many tests fail due to spurious reasons.

**Generating failure-inducing rules.** Grammar-based test generation [52] has been shown to be effective for avoiding spurious failures in fuzz testing. More recently, grammars and probabilistic variations of grammars have been used to infer abstract failure-inducing rules [66]. These rules can assist with the diagnosis of system failures, serve as accurate and high-level test oracles, and enable programmers to validate their fixes and prevent overfitting [51, 65, 66]. Our work takes inspiration from the research on inferring failure-inducing rules, but differs from the existing work on this topic in important ways. First, we focus on systems with numeric inputs, whereas existing research primarily deals with string-based inputs governed by a grammar. Second, instead of relying on an input grammar to generate tests, we investigate various test-generation heuristics that are guided by quantitative fitness functions drawn from system requirements. An exception is the work of Böhme et al. [36], which infers program patches for numeric systems without the need for input grammars. However, this approach relies on the availability of a human oracle to validate the verdicts of individual test inputs. In our context, this would be expensive and likely infeasible. Our work further differs from the above in that our goal is to identify rules for spurious failures rather than generating program patches.

The closest work to ours is the Alhazen framework [65], which we compared with in RQ3. In addition to the discussion and empirical comparison in RQ3, we note that our approach differs from Alhazen in its input-feature engineering for failure models. We derive the input features for decision-rule models from domain-knowledge heuristics, whereas Alhazen derives the input features dynamically from its input grammar. While Alhazen automates input-feature engineering, by incorporating domain knowledge into the design of input features, our approach provides the flexibility to derive rules that more closely match expert intuition.

**Test Input Validation.** Test input validation determines whether the test inputs given to a system adhere to the format, range or constraints anticipated by the system requirements. Test input validation improves the reliability of test results and is crucial for software testing [81]. Traditional software testing and verification approaches for CPS and network systems assume that pre-conditions describing valid inputs are already specified [28, 73, 74] or, alternatively, rely

on formal assume-guarantee and design-by-contract techniques [42, 56, 82]. Techniques based on assume-guarantee and design by contract require high-level formal system specifications. Such specifications do not exist for many real-world systems including our study subjects. Recent studies explore test input validity for deep learning (DL) models [59, 81], demonstrating that existing DL testing techniques generate several invalid test inputs. To mitigate this problem, the studies preform human subject experiments and establish metrics that determine test input validity for DL models. Although our work is not concerned with DL models, our definition of spurious test inputs is similar to that of invalid inputs for DL models [59, 81]. Similar to DL testing, failing to account for input validity leads to the generation of many invalid test inputs, thus reducing the reliability of test generation. In addition, similar to the research for DL models, we identify the rules leading to spurious failures for our study subjects based on domain expertise and human knowledge.

## 7   CONCLUSION

In this paper, we presented a data-driven framework for inferring failure models for systems with numeric inputs including cyber-physical and network systems. The framework employs existing surrogate-assisted and machine learning-guided (ML-guided) test generation techniques. We proposed a new dynamic surrogate-assisted algorithm that uses multiple surrogate models simultaneously during search, and dynamically selects the predictions from the most accurate model. We compared the accuracy of failure models obtained using our dynamic surrogate-assisted approach against two ML-guided techniques as well as two baselines using 16 study subjects from the cyber-physical and network domains. Our results, confirmed by statistical tests, show that the average accuracy, precision and recall of the dynamic surrogate-assisted approach are higher than those of the ML-guided test generation algorithms, and of the state-of-the-art and random-search baselines. Moreover, the rules inferred from the failure models built for our compute-intensive subjects identify genuine spurious failures as validated against domain knowledge. For future work, we plan to apply our approach to computer-vision and autonomous-driving systems and subsequently use the rules inferred by failure models as guidance for generating test inputs.

## 8   DATA AVAILABILITY

Implementations of all algorithms are available at [18]. The Simulink benchmark is available at [2], and NTSS at [7]. The requirements specifications for all study subjects are provided at [6]. Our evaluation data includes: (1) raw datasets for the experiments [13]; (2) rules generated for CI subjects [17, 20]; (3) evaluation scripts [3] and the analyzed data [15]; (4) statistical analysis results [16, 23, 26]; and (5) scripts for the plots in the paper [3].

## REFERENCES

[1] (Accessed: June 2023). *Autopilot online benchmark*. https://www.mathworks.com/matlabcentral/fileexchange/41490-autopilot-demo-for-arp4754a-do-178c-and-do-331?focused=6796756&tab=model

[2] (Accessed: June 2023). *Benchmark for Simulink models*. https://github.com/anonpaper23/testGenStrat/tree/main/Benchmark/Simulink%20Models

[3] (Accessed: June 2023). *Code to generate results of each research questions*. https://github.com/anonpaper23/testGenStrat/tree/main/Evaluation

[4] (Accessed: June 2023). *Code to SoTA implementation for NTSS case study*. https://github.com/anonpaper23/testGenStrat/blob/main/Code/NTSS/SoTA.py

[5] (Accessed: June 2023). *Code to SoTA implementation for Simulink model case study.* https://github.com/anonpaper23/testGenStrat/blob/main/Code/Simulink/Algorithms/decisiontreeSoTA.m

[6] (Accessed: June 2023). *CPS and NTSS requirements.* https://github.com/anonpaper23/testGenStrat/blob/main/Benchmark/Formalization/CPS_and_NTSS_Formalization.pdf

[7] (Accessed: June 2023). *ENRICH – non-robustnEss aNalysis for tRaffIC sHaping.* https://github.com/baharin/ENRICH

[8] (Accessed: June 2023). *Figure 16 to Figure 21 – precision and recall results obtained by varying time budget in RQ2.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[9] (Accessed: June 2023). *Figure 9 – Comparing Dataset sizes for dynamic SA algorithm and seven individual SA algorithms in RQ1.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[10] (Accessed: June 2023). *Lockheed Martin.* https://www.lockheedmartin.com

[11] (Accessed: June 2023). *Logistic Regression.* http://faculty.cas.usf.edu/mbrannick/regression/Logistic.html

[12] (Accessed: June 2023). *OpenWrt.* www.openwrt.org

[13] (Accessed: June 2023). *Raw datasets obtained from each algorithm for CPS and NTSS.* https://github.com/anonpaper23/testGenStrat/tree/main/Data/Dataset

[14] (Accessed: June 2023). *Replication package of Alhazen framework.* https://zenodo.org/records/3902142

[15] (Accessed: June 2023). *Results of each research question.* https://github.com/anonpaper23/testGenStrat/tree/main/Evaluation%20Results

[16] (Accessed: June 2023). *Results of statistical analysis.* https://github.com/anonpaper23/testGenStrat/blob/main/Evaluation%20Results/RQ2/RQ2StatisticalResults.xlsx

[17] (Accessed: June 2023). *Rules obtained for each CI subject.* https://github.com/anonpaper23/testGenStrat/blob/main/Evaluation%20Results/RQ4/APandNTSS_Rules.xlsx

[18] (Accessed: June 2023). *Source codes of algorithms for CPS and NTSS.* https://github.com/anonpaper23/testGenStrat/tree/main/Code

[19] (Accessed: June 2023). *Table 15 to Table 20 – average accuracy, recall and precision over all runs of algorithms by varying execution time budget in RQ2.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[20] (Accessed: June 2023). *Table 21 to Table 24 – full set of rules obtained for NTSS, AP1, AP2 and AP3 in RQ4.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[21] (Accessed: June 2023). *Table 3 – Parameter names, descriptions and values used by SoTA.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[22] (Accessed: June 2023). *Table 5 – Time budgets given to non-CI subjects in RQ1.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[23] (Accessed: June 2023). *Table 6 – Statistical tests for dataset size and percentage of incorrect labels over dataset size in RQ1.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[24] (Accessed: June 2023). *Table 7 – Time budget considered for CI subjects in RQ2.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[25] (Accessed: June 2023). *Table 8 – maximum number of test executions for non-CI subjects in RQ2.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[26] (Accessed: June 2023). *Table 9 to Table 14 – statistical tests for accuracy, recall and precision by varying execution time budget in RQ2.* https://github.com/anonpaper23/testGenStrat/blob/main/Supplementary_Material.pdf

[27] (Accessed: June 2023). tc-cake. https://man7.org/linux/man-pages/man8/tc-cake.8.html

[28] Yashwanth Annpureddy, Che Liu, Georgios Fainekos, and Sriram Sankaranarayanan. 2011. S-TaLiRo: A Tool for Temporal Logic Falsification for Hybrid Systems. In *Tools and Algorithms for the Construction and Analysis of Systems*, Parosh Aziz Abdulla and K. Rustan M. Leino (Eds.). Springer, 254–257.

[29] Aitor Arrieta, Shuai Wang, Urtzi Markiegi, Ainhoa Arruabarrena, Leire Etxeberria, and Goiuria Sagardui. 2019. Pareto efficient multi-objective black-box test case selection for simulation-based testing. *Information and Software Technology* (2019).

[30] Aitor Arrieta, Shuai Wang, Urtzi Markiegi, Goiuria Sagardui, and Leire Etxeberria. 2017. Search-based test case generation for cyber-physical systems. In *2017 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 688–697.

[31] Federal Aviation Administration (FAA)/Aviation Supplies & Academics (ASA). 2009. *Advanced Avionics Handbook*. Aviation Supplies & Academics, Incorporated. https://books.google.lu/books?id=2xGuPwAACAAJ

[32] Cornelius Aschermann, Tommaso Frassetto, Thorsten Holz, Patrick Jauernig, Ahmad-Reza Sadeghi, and Daniel Teuchert. 2019. NAUTILUS: Fishing for Deep Bugs with Grammars.. In *NDSS*.

[33] Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. 2017. Synthesizing program input grammars. *ACM SIGPLAN Notices* 52, 6 (2017), 95–110.

[34] Halil Beglerovic, Michael Stolz, and Martin Horn. 2017. Testing of autonomous vehicles using surrogate models and stochastic optimization. In *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 1–6.

[35] Raja Ben Abdessalem, Shiva Nejati, Lionel C Briand, and Thomas Stifter. 2016. Testing advanced driver assistance systems using multi-objective search and neural networks. In *Proceedings of the 31st IEEE/ACM international conference on automated software engineering*. 63–74.

[36] Marcel Böhme, Charaka Geethal, and Van-Thuan Pham. 2020. Human-in-the-loop automatic program repair. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*. IEEE, 274–285.

[37] Caius Brindescu, Iftekhar Ahmed, Rafael Leano, and Anita Sarma. 2020. Planning for untangling: Predicting the difficulty of merge conflicts. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 801–811.

[38] Devendra K Chaturvedi. 2017. *Modeling and simulation of systems using MATLAB® and Simulink®*. CRC press.

[39] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research* 16 (2002), 321–357.

[40] William W Cohen. 1995. Fast effective rule induction. In *Machine learning proceedings 1995*. Elsevier, 115–123.

[41] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 337–340.

[42] Patricia Derler, Edward A Lee, Stavros Tripakis, and Martin Törngren. 2013. Cyber-physical system design contracts. In *Proceedings of the ACM/IEEE 4th International Conference on Cyber-Physical Systems*. 109–118.

[43] Alan Díaz-Manríquez, Gregorio Toscano, Jose Hugo Barron-Zambrano, and Edgar Tello-Leal. 2016. A review of surrogate assisted multiobjective evolutionary algorithms. *Computational intelligence and neuroscience* 2016 (2016).

[44] Arkadiy Dushatskiy, Tanja Alderliesten, and Peter AN Bosman. 2021. A novel surrogate-assisted evolutionary algorithm applied to partition-based ensemble learning. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 583–591.

[45] Robert Feldt and Shin Yoo. 2020. Flexible Probabilistic Modeling for Search Based Test Data Generation. In *Proceedings of the 13th International Workshop on Search-Based Software Testing (SBST)*. 537–540.

[46] Martina Friese, Thomas Bartz-Beielstein, and Michael Emmerich. 2016. Building ensembles of surrogates by optimal convex combination. *Bioinspired optimization methods and their applications* (2016), 131–143.

[47] Khouloud Gaaloul, Claudio Menghi, Shiva Nejati, Lionel C Briand, and David Wolfe. 2020. Mining assumptions for software components using machine learning. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 159–171.

[48] Baljinder Ghotra, Shane McIntosh, and Ahmed E Hassan. 2015. Revisiting the impact of classification techniques on the performance of defect prediction models. In *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, Vol. 1. IEEE, 789–800.

[49] Dimitra Giannakopoulou, Corina S Pasareanu, and Howard Barringer. 2002. Assumption generation for software component verification. In *International Conference on Automated Software Engineering*. IEEE, 3–12.

[50] Dimitra Giannakopoulou, Thomas Pressburger, Anastasia Mavridou, and Johann Schumann. 2021. Automated formalization of structured natural language requirements. *Information and Software Technology* 137 (2021), 106590. https://doi.org/10.1016/j.infsof.2021.106590

[51] Rahul Gopinath, Alexander Kampmann, Nikolas Havrikov, Ezekiel O Soremekun, and Andreas Zeller. 2020. Abstracting failure-inducing inputs. In *Proceedings of the 29th ACM SIGSOFT international symposium on software testing and analysis*. 237–248.

[52] Kenneth V. Hanford. 1970. Automatic generation of test cases. *IBM Systems Journal* 9, 4 (1970), 242–257.

[53] Fitash Ul Haq, Donghwan Shin, Shiva Nejati, and Lionel Briand. 2021. Can Offline Testing of Deep Neural Networks Replace Their Online Testing? A Case Study of Automated Driving Systems. *Empirical Software Engineering* 26, 5 (2021), 90.

[54] Mark Harman, Sung Gon Kim, Kiran Lakhotia, Phil McMinn, and Shin Yoo. 2010. Optimizing for the number of tests generated in search based test data generation with an application to the oracle cost problem. In *2010 Third International Conference on Software Testing, Verification, and Validation Workshops*. IEEE, 182–191.

[55] Mark Harman and Phil McMinn. 2009. A theoretical and empirical study of search-based testing: Local, global, and hybrid search. *IEEE Transactions on Software Engineering* 36, 2 (2009), 226–247.

[56] Thomas A Henzinger, Shaz Qadeer, and Sriram K Rajamani. 1998. You assume, we guarantee: Methodology and case studies. In *Computer Aided Verification: 10th International Conference, CAV'98 Vancouver, BC, Canada, June 28–July 2, 1998 Proceedings 10*. Springer, 440–451.

[57] Toke Høiland-Jørgensen, Dave Täht, and Jonathan Morton. 2018. Piece of CAKE: a comprehensive queue management solution for home gateways. In *2018 IEEE International Symposium on Local and Metropolitan Area Networks (LANMAN)*. IEEE, 37–42.

[58] Linxiong Hong, Huacong Li, and Jiangfeng Fu. 2022. A novel surrogate-model based active learning method for structural reliability analysis. *Computer Methods in Applied Mechanics and Engineering* 394 (2022), 114835. https://doi.org/10.1016/j.cma.2022.114835

[59] Boyue Caroline Hu, Lina Marsso, Krzysztof Czarnecki, Rick Salay, Huakun Shen, and Marsha Chechik. 2022. If a Human Can See It, so Should Your System: Reliability Requirements for Machine Vision Components. In *Proceedings of the 44th International Conference on Software Engineering* (Pittsburgh, Pennsylvania) *(ICSE '22)*. Association for Computing Machinery, New York, NY, USA, 1145–1156. https://doi.org/10.1145/3510003.3510109

[60] Dmytro Humeniuk, Giuliano Antoniol, and Foutse Khomh. 2021. Data driven testing of cyber physical systems. In *2021 IEEE/ACM 14th International Workshop on Search-Based Software Testing (SBST)*. IEEE, 16–19.

[61] Dmytro Humeniuk, Foutse Khomh, and Giuliano Antoniol. 2022. A search-based framework for automatic generation of testing environments for cyber-physical systems. *Information and Software Technology* (2022), 106936.

[62] Yaochu Jin. 2005. A comprehensive survey of fitness approximation in evolutionary computation. *Soft computing* 9, 1 (2005), 3–12.

[63] Yaochu Jin and Bernhard Sendhoff. 2002. Fitness Approximation In Evolutionary Computation-a Survey.. In *GECCO*, Vol. 2. 1105–12.

[64] Baharin A. Jodat, Shiva Nejati, Mehrdad Sabetzadeh, and Patricio Saavedra. 2023. Learning Non-robustness using Simulation-based Testing: a Network Traffic-shaping Case Study. In *2023 IEEE Conference on Software Testing, Verification and Validation (ICST)*. IEEE, 386–397.

[65] Alexander Kampmann, Nikolas Havrikov, Ezekiel O Soremekun, and Andreas Zeller. 2020. When does my program do this? learning circumstances of software behavior. In *Proceedings of the 28th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering*. 1228–1239.

[66] Charaka Geethal Kapugama, Van-Thuan Pham, Aldeida Aleti, and Marcel Böhme. 2022. Human-in-the-loop oracle learning for semantic bugs in string processing programs. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*. 215–226.

[67] Fitsum Meshesha Kifetew, Roberto Tiella, and Paolo Tonella. 2017. Generating valid grammar-based test inputs by means of genetic programming and annotated grammars. *Empirical Software Engineering* 22, 2 (2017), 928–961.

[68] Neil Kulkarni, Caroline Lemieux, and Koushik Sen. 2021. Learning Highly Recursive Input Grammars. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 456–467.

[69] Jaekwon Lee, Seung Yeob Shin, Shiva Nejati, Lionel C Briand, and Yago Isasi Parache. 2022. Estimating Probabilistic Safe WCET Ranges of Real-Time Systems at Design Stages. *ACM Transactions on Software Engineering and Methodology* (2022).

[70] Sean Luke. 2013. *Essentials of Metaheuristics* (second ed.). Lulu. http://cs.gmu.edu/~sean/book/metaheuristics/.

[71] Reza Matinnejad, Shiva Nejati, and Lionel C. Briand. 2017. Automated testing of hybrid Simulink/Stateflow controllers: industrial case studies. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2017, Paderborn, Germany, September 4-8, 2017*. ACM, 938–943.

[72] Patrick E McKnight and Julius Najab. 2010. Mann-Whitney U Test. *The Corsini encyclopedia of psychology* (2010), 1–1.

[73] Claudio Menghi, Shiva Nejati, Lionel Briand, and Yago Isasi Parache. 2020. Approximation-refinement testing of compute-intensive cyber-physical models: An approach based on system identification. In *2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE)*. IEEE, 372–384.

[74] Claudio Menghi, Shiva Nejati, Khouloud Gaaloul, and Lionel C Briand. 2019. Generating automated and online test oracles for simulink models with continuous and uncertain behaviors. In *Proceedings of the 2019 27th acm joint meeting on european software engineering conference and symposium on the foundations of software engineering*. 27–38.

[75] Barton P Miller, Lars Fredriksen, and Bryan So. 1990. An empirical study of the reliability of UNIX utilities. *Commun. ACM* 33, 12 (1990), 32–44.

[76] Christoph Molnar. 2020. *Interpretable machine learning*. Lulu. com.

[77] Shiva Nejati, Khouloud Gaaloul, Claudio Menghi, Lionel C Briand, Stephen Foster, and David Wolfe. 2019. Evaluating model testing and model checking for finding requirements violations in Simulink models. In *Proceedings of the 2019 27th acm joint meeting on european software engineering conference and symposium on the foundations of software engineering*. 1015–1025.

[78] Shiva Nejati, Lev Sorokin, Damir Safin, Federico Formica, Mohammad Mahdi Mahboob, and Claudio Menghi. 2023. Reflections on Surrogate-Assisted Search-Based Testing: A Taxonomy and Two Replication Studies based on Industrial ADAS and Simulink Models. *Inf. Softw. Technol.* 163 (2023), 107286.

[79] Andrew Ng. 2018. Machine learning yearning. *Available: http://www.mlyearning.org/* (2018).

[80] Ripon Patgiri, Hemanth Katari, Ronit Kumar, and Dheeraj Sharma. 2019. Empirical study on malicious URL detection using machine learning. In *International Conference on Distributed Computing and Internet Technology*. Springer, 380–388.

[81] Vincenzo Riccio and Paolo Tonella. 2022. When and Why Test Generators for Deep Learning Produce Invalid Inputs: an Empirical Study. arXiv:2212.11368 [cs.SE]

[82] Alberto Sangiovanni-Vincentelli, Werner Damm, and Roberto Passerone. 2012. Taming Dr. Frankenstein: Contract-based design for cyber-physical systems. *European journal of control* 18, 3 (2012), 217–238.

[83] Alexander Schaap, Gordon Marks, Vera Pantelic, Mark Lawford, Gehan Selim, Alan Wassyng, and Lucian Patcas. 2018. Documenting Simulink Designs of Embedded Systems. In *International Conference on Model Driven Engineering Languages and Systems (MODELS): Companion Proceedings*. ACM, 47–51.

[84] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems* 25 (2012).

[85] Robert C Streijl, Stefan Winkler, and David S Hands. 2016. Mean opinion score (MOS) revisited: methods and applications, limitations and alternatives. *Multimedia Systems* 22, 2 (2016), 213–227.

[86] Hao Tong, Changwu Huang, Leandro L Minku, and Xin Yao. 2021. Surrogate models in evolutionary single-objective optimization: A new taxonomy and experimental study. *Information Sciences* 562 (2021), 414–437.

[87] Cumhur Erkan Tuncali, Georgios Fainekos, Danil Prokhorov, Hisahiro Ito, and James Kapinski. 2019. Requirements-driven test generation for autonomous vehicles with machine learning components. *IEEE Transactions on Intelligent Vehicles* 5, 2 (2019), 265–280.

[88] András Vargha and Harold D Delaney. 2000. A critique and improvement of the CL common language effect size statistics of McGraw and Wong. *Journal of Educational and Behavioral Statistics* 25, 2 (2000), 101–132.

[89] Junjie Wang, Bihuan Chen, Lei Wei, and Yang Liu. 2019. Superion: Grammar-aware greybox fuzzing. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 724–735.

[90] Yan Wang, Peng Jia, Luping Liu, Cheng Huang, and Zhonglin Liu. 2020. A systematic review of fuzzing based on machine learning techniques. *PloS one* 15, 8 (2020), e0237749.

[91] Ian H. Witten, Eibe Frank, and Mark A. Hall. 2011. *Data Mining: Practical Machine Learning Tools and Techniques* (3 ed.). Morgan Kaufmann, Amsterdam. http://www.sciencedirect.com/science/book/9780123748560

[92] Huanwei Xu, Xin Zhang, Hao Li, and Ge Xiang. 2021. An ensemble of adaptive surrogate models based on local error expectations. *Mathematical Problems in Engineering* 2021 (2021).