

AIM: Automated Input Set Minimization for Metamorphic Security Testing

Nazanin Bayati Chaleshtari, Yoann Marquer, Fabrizio Pastore, *Member, IEEE*, and Lionel C. Briand, *Fellow, IEEE*

Abstract—Although the security testing of Web systems can be automated by generating crafted inputs, solutions to automate the test oracle, i.e., vulnerability detection, remain difficult to apply in practice. Specifically, though previous work has demonstrated the potential of metamorphic testing—security failures can be determined by metamorphic relations that turn valid inputs into malicious inputs—metamorphic relations are typically executed on a large set of inputs, which is time-consuming and thus makes metamorphic testing impractical.

We propose AIM, an approach that automatically selects inputs to reduce testing costs while preserving vulnerability detection capabilities. AIM includes a clustering-based black-box approach, to identify similar inputs based on their security properties. It also relies on a novel genetic algorithm to efficiently select diverse inputs while minimizing their total cost. Further, it contains a problem-reduction component to reduce the search space and speed up the minimization process. We evaluated the effectiveness of AIM on two well-known Web systems, Jenkins and Joomla, with documented vulnerabilities. We compared AIM's results with four baselines involving standard search approaches. Overall, AIM reduced metamorphic testing time by 84% for Jenkins and 82% for Joomla, while preserving the same level of vulnerability detection. Furthermore, AIM significantly outperformed all the considered baselines regarding vulnerability coverage.

Index Terms—System Security Testing, Metamorphic Testing, Test Suite Minimization, Many-Objective Search

I. INTRODUCTION

Web systems, from social media platforms to e-commerce and banking systems, are a backbone of our society: they manage data that is at the heart of our social and business activities (e.g., public pictures, bank transactions), and, as such, should be protected. To verify that Web systems are secure, engineers perform security testing, which consists of verifying that the software adheres to its security properties (e.g., confidentiality, availability, and integrity). Such testing is typically performed by simulating malicious users interacting with the system under test [1], [2].

At a high-level, security testing does not differ from other software testing activities: it consists of providing inputs to the software under test and verifying that the software outputs

are correct, based on specifications. For such a verification, a *test oracle* [3] is required, i.e., a mechanism for determining whether a test case has passed or failed. When test cases are manually derived, test oracles are defined by engineers and they generally consist of expressions comparing an observed output with the expected output, determined from software specifications. In security testing, when a software output differs from the expected one, then a software vulnerability (i.e., a fault affecting the security properties of the software) has been discovered.

Deriving test oracles for the software under test (SUT) is called the *oracle problem* [4], which entails distinguishing correct from incorrect outputs for all potential inputs. Except for the verification of basic reliability properties—ensuring that the software provides a timely output and does not crash—the problem is not tractable without additional executable specifications (e.g., method post-conditions or detailed system models), which, unfortunately, are often unavailable. Further, since software vulnerabilities tend to be subtle, it is necessary to exercise each software interface with a large number of inputs (e.g., providing all the possible code injection strings to a Web form). When a large number of test inputs are needed, even in the presence of automated means to generate them (e.g., catalogs of code injection strings), testing becomes impractical if we lack solutions to automatically derive test oracles.

Metamorphic testing was proposed to alleviate the oracle problem [5] by testing not the input-output behavior of the system, but by comparing the outputs of multiple test executions [5], [6]. It relies on metamorphic relations (MRs), which are specifications expressing how to derive a follow-up input from a source input and relations between the corresponding outputs. Such an approach has shown to be useful for security testing, also referred to as metamorphic security testing (MST) [6], [7]. MST consists in relying on MRs to modify source inputs to obtain follow-up inputs that mimic attacks and verify that known output properties captured by these MRs hold (e.g., if the follow-up input differs from the source input in some way, then the output shall be different). For instance, one may verify if URLs can be accessed by users who should not reach them through their user interface, thus enabling the detection of authorization vulnerabilities.

MST has been successfully applied to testing Web interfaces [6], [7] in an approach called MST-wi [6]; in such context, source inputs are sequences of interactions with a Web system and can be easily derived using a Web crawler. For example, a source input may consist of two actions: performing a login and then requesting a specific URL appearing

N. Bayati Chaleshtari and L. Briand are with the School of Electrical and Computer Engineering of University of Ottawa, Canada, Y. Marquer and F. Pastore are with the Interdisciplinary Centre for Security, Reliability, and Trust (SnT) of the University of Luxembourg, Luxembourg, and L. Briand is also with Lero SFI Centre for Software Research and University of Limerick, Ireland. Part of this work was done when L. Briand was affiliated with the Interdisciplinary Centre for Security, Reliability, and Trust (SnT) of the University of Luxembourg.
E-mail: n.bayati@uottawa.ca, yoann.marquer@uni.lu, fabrizio.pastore@uni.lu, lbriand@uottawa.ca

Manuscript received Month DD, 2024; revised Month DD, 2024.

in the returned Web page. MST-wi integrates a catalog of 76 MRs enabling the identification of 101 vulnerability types.

The performance and scalability of metamorphic testing naturally depends on the number of source inputs to be processed. In the case of MST-wi, we demonstrated that scalability can be achieved through parallelism; however, such solution may not fit all development contexts (e.g., not all companies have an infrastructure enabling parallel execution of the software under test and its test cases). Further, even when parallelization is possible, a reduction of the test execution time may provide tangible benefits, including earlier vulnerability detection. In general, what is required is an approach to minimize the number of source inputs to be used during testing.

In this work, we address the problem of minimizing source inputs used by MRs to make metamorphic testing more scalable, with a focus on Web systems, though many aspects are reusable to other domains. We propose the Automated Input Minimizer (AIM) approach, which aims at minimizing a set of source inputs (hereafter, the *initial input set*), while preserving the capability of MRs to detect security vulnerabilities. More in detail, this work includes the following contributions:

- We propose AIM, an approach to minimize input sets for metamorphic testing while retaining inputs able to detect vulnerabilities. Note that many steps of AIM are not specific to Web systems while others would need to be tailored to other domains (e.g., desktop applications, embedded systems). This approach includes the following novel components:
 - An extension of the MST-wi framework to retrieve output data and extract cost information about MRs without executing them.
 - A black-box approach leveraging clustering algorithms to partition the initial input set based on security-related characteristics in order to keep a small number of representative inputs.
 - MOCCO (Many-Objective Coverage and Cost Optimizer), a novel genetic algorithm specifically designed to fit our problem and which is able to efficiently select diverse inputs while minimizing their total cost.
 - IMPRO (Input set Minimization Problem Reduction Operator), an approach to reduce the search space to its minimal extent, and then divide it in smaller, easier to minimize, independent parts.
- We provide a prototype framework for AIM [8], integrating the above components and automating the process of input set minimization for Web systems.
- We report on an extensive empirical evaluation (about 800 hours of computation) aimed at assessing the effectiveness of AIM in terms of vulnerability detection and performance, considering 18 different AIM configurations and 5 search algorithms (including MOCCO) for security testing, on the Jenkins and Joomla systems, which are the most used Web-based frameworks for development automation and context management.
- We also provide a proof of the correctness of the AIM approach in a separate appendix provided as supplementary

material.

This paper is structured as follows. We introduce background information necessary to state our problem and detail our approach (Section II). We define the problem of minimizing the initial input set while retaining inputs capable of detecting distinct software vulnerabilities (Section III). We present an overview of AIM (Section IV) and then detail our core technical solutions (Sections V to IX). We report on a large-scale empirical evaluation of AIM (Section X) and address the threats to the validity of the results (Section XI). We discuss and contrast related work (Section XII) and draw conclusions (Section XIII).

II. BACKGROUND

In this section, we present the concepts required to define our approach. We first provide a background on Metamorphic Testing (MT, § II-A), then we briefly describe MST-wi, our previous work on the application of MT to security (§ II-B). Next, we briefly describe three clustering algorithms: K-means (§ II-D1), DBSCAN (§ II-D2), and HDBSCAN (§ II-D3). Finally, we introduce optimization problems (§ II-E).

A. Metamorphic Testing

In contrast to common testing practice, which compares for each input of the system the actual output against the expected output, MT examines the relationships between outputs obtained from multiple test executions.

MT is based on Metamorphic Relations (MRs), which are necessary properties of the SUT (system under test) in relation to multiple inputs and their expected outputs [9]. The test result, either pass or failure, is determined by validating the outputs of various executions against the MR.

Formally, let S be the SUT. In the context of MT, inputs in the domain of S are called *source inputs*. Moreover, we call *source output* and we denote $S(x)$ the output obtained from a source input x . An MR is the combination of:

- A *transformation function* θ , taking values in source inputs and generating new inputs called *follow-up inputs*. For each source input x , we call *follow-up output* the output $S(\theta(x))$ of the follow-up input $\theta(x)$.
- An *output relation* R between source outputs and follow-up outputs.

The MR is *executed* with a source input x when the follow-up input $\theta(x)$ is generated, then the SUT is executed on both inputs to obtain outputs $S(x)$ and $S(\theta(x))$, and finally the relation $R(S(x), S(\theta(x)))$ is checked. If this relation holds, then the MR is *satisfied*, otherwise it is *violated*.

For instance, consider a system implementing the cosine function. It might not be feasible to verify the $\cos(x)$ results for all possible values of x , except for special values of x , e.g., $\cos(0) = 1$ or $\cos(\frac{\pi}{2}) = 0$. However, the cosine function satisfies that, for each input x , $\cos(\pi - x) = -\cos(x)$. Based on this property, we can define an MR, where the source inputs are the possible angle values of x , the follow-up inputs are $y = \pi - x$, and the expected relation between source and follow-up outputs is $\cos(y) = -\cos(x)$. The SUT is executed

```

1 MR CWE_668 {
2 {
3   var sep = "/";
4   for (var par=0; par < 4; par++){
5     for (Action action : Input(1).actions()){
6       var pos = action.getPosition();
7       var newUrl = action.urlPath+sep+RandomFilePath();
8       IMPLIES(
9         !isAdmin(action.user) &&
10        afterLogin(action) &&
11        CREATE(Input(2), Input(1)) &&
12        Input(2).actions().get(pos).setUrl(newUrl) &&
13        notTried(action.getUser(), newUrl)
14      ,
15      TRUE(
16        Output(Input(2),pos).noFile() ||
17        userCanRetrieveContent(action.getUser(), Output(Input(2),pos).file()) ||
18        different(Output(Input(1),pos), Output(Input(2),pos))
19      );//end-IMPLIES
20    }//end-for
21    sep=sep+"../";
22  }//end-for
23 }//end-MR
24 }//end-package

```

Fig. 1. MR for CWE 668: Exposure of resource to wrong sphere [12]

twice per source input, respectively with an angle x and an angle $y = \pi - x$. The outputs of both executions are then validated against the output relation. If this relation is violated, then the SUT is faulty.

B. Metamorphic Security Testing

In previous work, we automated MT in the security domain by introducing a tool named MST-wi [6]. MST-wi enables software engineers to define MRs that capture the security properties of Web systems. MST-wi includes a data collection framework that crawls the Web system under test to automatically derive source inputs. Each source input is a sequence of interactions of the legitimate user with the Web system. Moreover, MST-wi includes a Domain Specific Language (DSL) to support writing MRs for security testing. Finally, MST-wi provides a testing framework that automatically performs security testing based on the defined MRs and the input data.

In MST, follow-up inputs are generated by modifying source inputs, simulating actions an attacker might take to identify vulnerabilities in the SUT. These modifications can be done using 55 Web-specific functions enabling engineers to define complex security properties, e.g., `cannotReachThroughGUI`, `isSupervisorOf`, and `isError`. MRs capture security properties that hold when the SUT behaves in a safe way. If an MR, for any given source input, gets violated, then MST-wi detects a vulnerability in the SUT. In that case, we say that the MR *exercised* the vulnerability in the SUT. MST-wi includes a catalog of 76 MRs, inspired by OWASP guidelines [10] and vulnerability descriptions in the CWE database [11], capturing a large variety of security properties for Web systems.

We describe in Figure 1 an MR written for CWE 668, which concerns unintended access rights [12]. This MR verifies that a file path passed in a URL should never enable a user to access data that is not already provided by the user

interface. The first `for` loop iterates multiple times (Line 4) to cover different system paths, e.g., `/` and `/../..`. The second `for` loop iterates over all the actions of a source input (Line 5). Each action in the sequence is identified by its position, i.e., the first action in a sequence has position 0. The position of the current action is stored (Line 6) to be used to generate the corresponding follow-up action. A new URL is defined by concatenating the URL of the current action and a randomly selected system file path, e.g., `config.xml`, (Line 7). For instance, if the URL of the current action is `http://www.hostname.com`, the new URL can be `http://www.hostname.com/../../config.xml`.

The MR first checks if the user who is performing the action is admin (Line 9), since an admin has direct access to the system file path, and hence will not exercise a vulnerability. Then, the MR checks that the action is performed after a login (Line 10), to ensure this action requires authentication. Then, the MR generates a follow-up input, named `Input(2)`, by copying the current sequence of actions `Input(1)` (Line 11) and setting the URL of the current action to the new URL (Line 12). To speed up the process, the MR verifies that the current user has not tried the same URL before (Line 13). The SUT is vulnerable if all the following conditions are violated: 1) the follow-up input does not access a file at the new URL, or 2) it accesses a file, but the user has the right to access it, or 3) the source and follow-up inputs obtain different outputs, as the follow-up input tries to access a system file without access rights, while the source input is accessing the originally crawled URL.

This MR tests the initial set of source inputs, with different URLs and users, and transforms each one several times with different system file paths, leading to a combinatorial explosion. The more executed actions, the longer the execution time. The provided MR, with an input set of 160 source inputs on Jenkins, executed more than 200,000 follow-up

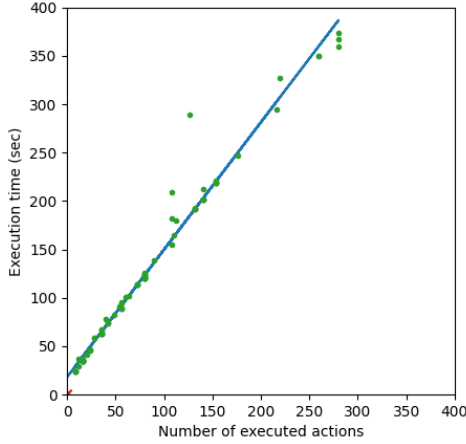


Fig. 2. Linear regression between the number of executed actions and the execution time of a metamorphic relation. Each input is represented by a green dot, while the blue line depicts the linear regression model.

inputs in 17,694 minutes (about 12 days) on a professional desktop PCs (Dell G7 7500, RAM 16Gb, Intel(R) Core(TM) i9-10885H CPU @ 2.40GHz). Even when parallelization is possible, a reduction of the test execution time may provide tangible benefits, including earlier vulnerability detection. This warrants an approach to minimize the initial set of source inputs, based on the cost of each input.

However, knowing the execution time of each source input would require to execute them on the considered MRs, hence defeating the purpose of input set minimization. To avoid executing MRs but relying on the number of actions executed by an MR as a surrogate metric for execution time to guide the minimization technique, we computed the Spearman's correlation coefficients between execution time and number of actions exercised by the MRs tested in a previous study [6]. It led to significant correlations (i.e., coefficients of correlation above 0.5 and p-value p below 0.05), thus confirming the feasibility of relying on the number of actions as a surrogate metric for execution time.

The example in Figure 2 depicts a typical linear correlation between the execution time of an MR and the number of executed actions. It uses randomly selected source inputs and an MR written for CWE 863. Each point represents the execution time (x-axis) and the number of executed actions (y-axis) for a given input. The linear regression is represented by the blue line and the corresponding coefficient of determination is 97.8%, indicating a strong linear correlation.

C. Test Suite Minimization

Test suites are prone to redundant test cases that, if not removed, can lead to a massive waste of time and resources [13], thus warranting systematic and automated strategies to eliminate redundant test cases, that are referred to as test suite minimization.

While test suite minimization techniques are very diverse [13], most of them are white-box approaches aiming at

minimizing the size of the test suite while maximizing code coverage. For instance, several test minimization approaches used greedy heuristics to select test cases based on their code coverage [14], [15]. Black-box approaches include the FAST-R family of scalable approaches that leverages a representation of test source code (or command line inputs) in a vector-space model [16] and the ATM approach that is based on the abstract syntax tree of test source code [17]. Both use similarity metrics to cluster and then select test cases.

In the context of Web systems, both the system source code and test code are not available to determine similarity between source inputs, warranting a different black-box approach able to cluster and select source inputs for these systems. Moreover, to make MST scalable, MR execution time should be minimized while preserving vulnerability detection, warranting an approach that minimizes the number of executed actions (§ II-B) while covering all the input clusters.

D. Clustering

Within the clustering steps in Section VI, we rely on three well-known clustering algorithms: K-means (§ II-D1), DBSCAN (§ II-D2), and HDBSCAN (§ II-D3).

1) *K-means*: K-means is a clustering algorithm which takes as input a set of data points and an integer k . K-means aims to assign data points to k clusters by maximizing the similarity between individual data points within each cluster and the center of the cluster, called centroid. The centroids are randomly initialized, then iteratively refined until a fixpoint is reached [18].

2) *DBSCAN*: DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is an algorithm that defines clusters using local density estimation. This algorithm takes as input a dataset and two configuration parameters: the distance threshold ϵ and the minimum number of neighbors n .

The distance threshold ϵ is used to determine the ϵ -neighborhood of each data point, i.e., the set of data points that are at most ϵ distant from it. There are three different types of data points in DBSCAN, based on the number of neighbors in the ϵ -neighborhood of a data point:

- Core** If a data point has a number of neighbors above n , it is then considered a core point.
- Border** If a data point has a number of neighbors below n , but has a core point in its neighborhood, it is then considered a border point.
- Noise** Any data point which is neither a core point nor a border point is considered noise.

A cluster consists of the set of core points and border points that can be reached through their ϵ -neighborhoods [19]. DBSCAN uses a single global ϵ value to determine the clusters. But, if the clusters have varying densities, this could lead to suboptimal partitioning of the data. HDBSCAN addresses this problem and we describe next.

3) *HDBSCAN*: HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) is an extension of DBSCAN (§ II-D2). As opposed to DBSCAN, HDBSCAN relies on different distance thresholds ϵ for each cluster, thus obtaining clusters of varying densities.

HDBSCAN first builds a hierarchy of clusters, based on various ϵ values selected in decreasing order. Then, based on such a hierarchy, HDBSCAN selects as final clusters the most persistent ones, where cluster persistence represents how long a cluster remains the same without splitting when decreasing the value of ϵ . In HDBSCAN, one has only to specify one parameter, which is the minimum number of individuals required to form a cluster, denoted by n [20]. Clusters with less than n individuals are considered noise and ignored.

E. Many Objective Optimization

Engineers are often faced with problems requiring to fulfill multiple objectives at the same time, called *multi-objective problems*. For instance, multi-objective search algorithms were used in test suite minimization approaches to balance cost, effectiveness, and other objectives [21], [22]. Multi-objective problems with at least (three or) four objectives are informally known as *many-objective problems* [23]. In both kind of problems, one needs a solution which is a good trade-off between the objectives. Hence, we first introduce the Pareto front of a decision space (§ II-E1). Then, we describe genetic algorithms able to solve many-objective problems (§ II-E2).

1) *Pareto Front*: Multi- and many-objective problems can be stated as minimizing several objective functions while taking values in a given decision space. The goal of multi-objective optimization is to approximate the Pareto Front in the objective space [23].

Formally, if D is the decision space and $f_1(\cdot), \dots, f_n(\cdot)$ are n objective functions defined on D , then the *fitness vector* of a decision vector $x \in D$ is $[f_1(x), \dots, f_n(x)]$, hereafter denoted $F(x)$. Moreover, a decision vector x_1 *Pareto-dominates* a decision vector x_2 (hereafter denoted $x_1 \succ x_2$) if 1) for each $1 \leq i \leq n$, we have $f_i(x_1) \leq f_i(x_2)$, and 2) there exists $1 \leq i \leq n$ such that $f_i(x_1) < f_i(x_2)$. If there exists no decision vector x_1 such that $x_1 \succ x_2$, we say that the decision vector x_2 is *non-dominated*. The *Pareto front* of D is the set $\{F(x_2) \mid x_2 \in D \text{ and } \forall x_1 \in D : x_1 \not\succ x_2\}$ of the fitness vectors of the non-dominated decision vectors. Finally, a *multi/many-objective problem* consists in:

$$\underset{x \in D}{\text{minimize}} F(x) = [f_1(x), \dots, f_n(x)]$$

where the minimize notation means that we want to find or at least approximate the non-dominated decision vectors, hence the ones having a fitness vector in the Pareto front [23].

2) *Solving Many-Objective Problems*: Multi-objective algorithms like NSGA-II [24] or SPEA2 [25], [26] are not effective in solving many-objective problems [27], [28] because of the following challenges:

- 1) The proportion of non-dominated solutions becomes exponentially large with an increased number of objectives. This reduces the chances of the search being stuck at a local optimum and may lead to a better convergence rate [23], but also slows down the search process considerably [27].
- 2) With an increased number of objectives, diversity operators (e.g., based on crowding distance or clustering) become computationally expensive [27].

- 3) If only a handful of solutions are to be found in a large-dimensional space, solutions are likely to be widely distant from each other. Hence, two distant parent solutions are likely to produce offspring solutions that are distant from them. In this situation, recombination operations may be inefficient and require crossover restriction or other schemes [27].

To tackle these challenges, several many-objective algorithms have been successfully applied within the software engineering community, like NSGA-III [27], [29] and MOSA [28].

NSGA-III [27], [29] is based on NSGA-II [24] and addresses these challenges by assuming a set of supplied or predefined reference points. Diversity (challenge 2) is ensured by starting the search in parallel from each of the reference points, assuming that largely spread starting points would lead to exploring all relevant parts of the Pareto front. For each parallel search, parents share the same starting point, so they are assumed to be close enough so that recombination operations (challenge 3) are more meaningful. Finally, instead of considering all solutions in the Pareto front, NSGA-III focuses on individuals which are the closest to the largest number of reference points. That way, NSGA-III considers only a small proportion of the Pareto front (addressing challenge 1).

Another many-objective algorithm, MOSA [28], does not aim to identify a single individual achieving a trade-off between objectives but a set of individuals, each satisfying one of the objectives. Such characteristic makes MOSA adequate for many software testing problems where it is sufficient to identify one test case (i.e., an individual) for each test objective (e.g., covering a specific branch or violating a safety requirement). To deal with challenge 1, MOSA relies on a preference criterion amongst individuals in the Pareto front, by focusing on 1) extreme individuals (i.e., test cases having one or more objective scores equal to zero), and 2) in case of tie, the shortest test cases. These best extreme individuals are stored in an archive during the search, and the archive obtained at the last generation is the final solution. Challenges 2 and 3 are addressed by focusing the search, on each generation, on the objectives not yet covered by individuals in the archive.

III. PROBLEM DEFINITION

As the time required to execute a set of considered MRs may be large (§ II-B), we aim to minimize the set of source inputs (hereafter, *input set*) to be used when applying MST to a Web system, given a set of MRs. In our context, each input is a sequence of actions used to communicate with the Web system and each action leads to a different Web page.

To ensure that a *minimized input set* can exercise the same vulnerabilities as the original one, intuitively, we should ensure that they belong to the same *input blocks*. Indeed, in software testing, after identifying an important characteristic to consider for the inputs, one can partition the input space in blocks, i.e., pairwise disjoint sets of inputs, such that inputs in the same block exercise the SUT in a similar way [30]. As the manual identification of relevant input blocks for a large system is extremely costly, we rely on clustering for that

purpose (Section VI). Since an input is a sequence of actions, it can exercise several input blocks. In the rest of the paper, we rely on the notion of *input coverage*, indicating the input blocks an input belongs to.

A. Assumptions and Goals

We assume we know, for each input in the initial input set, 1) its cost and 2) its coverage.

1) Because we want to make MST scalable, the cost $cost(in)$ of an input in corresponds to the execution time required to verify if the considered MRs are satisfied with this input. Because we aim to reduce this execution time without having to execute the MRs, as it would defeat the purpose of input set minimization, we use the number $nbActions(mr, in)$ of actions to be executed by an MR mr using input in as a surrogate metric for its execution time (see § II-B). We thus define the cost of an input as follows:

$$cost(in) \stackrel{\text{def}}{=} \sum_{mr \in MRs} nbActions(mr, in)$$

When $cost(in) = 0$, input in was not exercised by any MR due to the preconditions in these MRs. Hence, in is not useful for MST and can be removed without loss from the initial input set. Finally, the total cost of an input set I is $cost(I) \stackrel{\text{def}}{=} \sum_{in \in I} cost(in)$.

2) To minimize the cost of metamorphic testing, we remove unnecessary inputs from the initial input set, but we want to preserve all the inputs able to exercise distinct vulnerabilities. Hence, we consider, for each initial input in , its coverage $Cover(in)$. In our study, $Cover(in)$ is the set of input blocks the input in belongs to, and we determine these input blocks in Section VI using *double-clustering*. For now, we assume that the coverage of an input is known. The total coverage of an input set I is $Cover(I) \stackrel{\text{def}}{=} \cup_{in \in I} Cover(in)$.

We can now state our goals. We want to obtain a subset $I_{final} \subseteq I_{init}$ of the initial input set such that 1) I_{final} does not reduce total input coverage, i.e., $Cover(I_{final}) = Cover(I_{init})$ and 2) I_{final} has minimal cost, i.e., $cost(I_{final}) = \min\{cost(I) \mid I \subseteq I_{init} \wedge Cover(I) = Cover(I_{init})\}$. Note that a solution I_{final} may not be necessarily unique.

B. A Many-Objective Problem

To minimize the initial input set, we focus on the *selection* of inputs that belong to the same input blocks as the initial input set. A potential solution to our problem is an input set $I \subseteq I_{init}$. Obtaining a solution I able to reach full input coverage is straightforward since, for each block bl , one can simply select an input in $Inputs(bl) \stackrel{\text{def}}{=} \{in \in I_{init} \mid bl \in Cover(in)\}$. The hard part of our problem is to determine a combination of inputs able to reach full input coverage at a minimal cost. Hence, we have to consider an input set as a whole and not focus on individual inputs.

This is similar to the *whole suite* approach [31] targeting white-box testing. They use as objective the total number of covered branches. But, in our context, counting the number of uncovered blocks would consider as equivalent input sets that miss the same number of blocks, without taking into account

that it may be easier to cover some blocks than others (e.g., some blocks may be covered by many inputs, but some only by a few) or that a block may be covered by inputs with different costs. Thus, to obtain a combination of inputs that minimizes cost while preserving input coverage, we have to investigate how input sets cover each input block.

Hence, we are interested in covering each input block as an individual objective, in a way similar to the coverage of each code branch for white-box testing [28]. Because the total number of blocks to be covered is typically large (≥ 4), we deal with a *many-objective problem* [23]. This can be an advantage, because a many-objective reformulation of complex problems can reduce the probability of being trapped in local optima and may lead to a better convergence rate [28]. But this raises several challenges (§ II-E2) that we tackle while presenting our search algorithm (Section VIII).

C. Objective Functions

To provide effective guidance to a search algorithm, we need to quantify when an input set is closer to the objective of covering a particular block bl than another input set. In other words, if I_1 and I_2 are two input sets which do not cover bl but have the same cost, we need to determine which one is more desirable to achieve the goals introduced in § III-A by defining appropriate objective functions.

In general, adding an input to an input set would not only cover bl , but would also likely cover other blocks, that would then be covered by several inputs, thus introducing the possibility to remove some of them without affecting coverage. To track of how a given block bl is covered by inputs from a given input set I , we introduce the concept of *superposition* as $superpos(bl, I) \stackrel{\text{def}}{=} |Inputs(bl) \cap I|$. For instance, if $superpos(bl, I) = 1$, then there is only one input in I covering bl . In that case, this input is necessary to maintain the coverage of I . More generally, with the *redundancy* metric, we quantify how much an input in is necessary to ensure the coverage of an input set I it belongs to: $redundancy(in, I) \stackrel{\text{def}}{=} \min\{superpos(bl, I) \mid bl \in Cover(in)\} - 1$. The -1 is used to normalize the redundancy metric so that its range starts at 0. If $redundancy(in, I) = 0$, we say that in is *necessary* in I , otherwise we say that in is *redundant*. In the following, we denote $Redundant(I) \stackrel{\text{def}}{=} \{in \in I \mid redundancy(in, I) > 0\}$ the set of the redundant inputs in I .

To focus on the least costly input sets during the search (Section VIII), we quantify the gain obtained by removing redundant inputs. If I contains a redundant input in , then we call *removal step* a transition from I to $I \setminus \{in\}$. Otherwise, we say that I is already *reduced*. Unfortunately, given two redundant inputs in_1 and in_2 , removing in_1 may render in_2 necessary. Hence, when considering potential removal steps (e.g., removing either in_1 or in_2), one has to consider the order of these steps. We represent a *valid order* of removal steps by a sequence of inputs $[in_1, \dots, in_n]$ to be removed from I such that, for each $0 \leq i < n$, in_{i+1} is redundant in $I \setminus \{in_1, \dots, in_i\}$. We denote $ValidOrders(I)$ the set of valid orders of removal steps in I . Removing redundant inputs

in_1, \dots, in_n leads to a reduction of cost $cost(in_1) + \dots + cost(in_n)$. For each input set I , we consider the maximal *gain* from valid orders of removal steps:

$$gain(I) \stackrel{\text{def}}{=} \max \left\{ \sum_{1 \leq i \leq n} cost(in_i) \mid [in_1, \dots, in_n] \in ValidOrders(I) \right\}$$

To reduce the cost of computing this gain, we prove in the separate appendix (Theorem 1) that, to determine which orders of removal steps are valid, we can remove inputs in any arbitrary order, without having to resort to backtracking to previous inputs. Moreover, in our approach, we need to compute the gain only in situations when the number of redundant inputs is small (Section VIII), thus exhaustively computing the gain is tractable.

Adding an input $in_1 \in Inputs(bl)$ to I would lead to a gain but would also result in additional cost, hence warranting we consider the benefit-cost balance $gain(I \cup \{in_1\}) - cost(in_1)$ to evaluate how efficiently bl is covered by in_1 . More generally, we define the *potential* of I to efficiently cover bl as the maximum benefit-cost balance obtained by adding inputs in_1 in covering bl . But, as an input in_1 may be necessary to cover bl while leading to no or not enough removal steps, $gain(I \cup \{in_1\}) - cost(in_1)$ may be negative. To facilitate normalization, we need the potential to return a non-negative value, thus we shift all the benefit-cost balances for a given objective bl by adding a dedicated term. As the potential is a maximum, the worst case of $gain(I \cup \{in_2\}) - cost(in_2)$ is when $gain(I \cup \{in_2\}) = 0$ and $cost(in_2)$ is the minimal cost amongst the inputs able to cover bl . Hence, we obtain the following definition for the potential of I in covering bl :

$$potential(I, bl) \stackrel{\text{def}}{=} \max \{ gain(I \cup \{in_1\}) - cost(in_1) \mid in_1 \in Inputs(bl) \} + \min \{ cost(in_2) \mid in_2 \in Inputs(bl) \}$$

We thus use the potential to define the objective function associated with an objective bl .

To normalize our metrics, we rely on the normalization function $\omega(x) \stackrel{\text{def}}{=} \frac{x}{x+1}$, which is used to reduce a range of values from $[0, \infty)$ to $[0, 1)$ while preserving the ordering. When used during a search, it is less prone to precision errors and more likely to drive faster convergence towards an adequate solution than alternatives [32]. We use it to normalize the cost between 0 and 1 and the smaller is the normalized cost, the better a solution is. For coverage, since a high potential is more desirable, we use its complement $\frac{1}{x+1} = 1 - \omega(x)$ to reverse the order, so that the more potential an input set has, the lower its coverage objective function is. We thus define the objective function corresponding to a block bl_i as:

$$f_{bl_i}(I) \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } bl_i \in Cover(I) \\ \frac{1}{potential(I, bl_i)+1} & \text{otherwise} \end{cases}$$

The lower this value, the better. If bl_i is covered, then $f_{bl_i}(I) = 0$, otherwise $f_{bl_i}(I) > 0$. As expected, input sets that cover the objective are better than input sets that do not, and if two input sets do not cover the objective, then the normalized potential is used to break the tie.

D. Solutions to Our Problem

Each element in the decision space is an input set $I \subseteq I_{init}$, which is associated with a *fitness vector*:

$$F(I) \stackrel{\text{def}}{=} [\omega(cost(I)), f_{bl_1}(I), \dots, f_{bl_n}(I)]$$

where $Cover(I_{init}) = \{bl_1, \dots, bl_n\}$ denotes the n input blocks to be covered by input sets $I \subseteq I_{init}$.

Hence, we can define the Pareto front formed by the non-dominated solutions in our decision space (§ II-E1) and we formulate our problem definition as a many-objective optimization problem:

$$\underset{I \subseteq I_{init}}{\text{minimize}} F(I)$$

where the minimize notation means that we want to find or at least approximate the non-dominated decision vectors having a fitness vector on the Pareto front [23]. Because we want full input coverage (§ III-A), the ultimate goal is a non-dominated solution I_{final} such that:

$$F(I_{final}) = [\omega(cost_{\min}), 0, \dots, 0]$$

where $cost_{\min}$ is the cost of the cheapest subset of I_{init} with full input coverage.

IV. OVERVIEW OF THE APPROACH

As stated in our problem definition (Section III), we aim to reduce the cost of MST (§ II-B) by minimizing an initial input set without removing inputs that are required to exercise distinct security vulnerabilities. To do so, we need to tackle the following sub-problems (§ III-D):

- 1) For each initial input, we need to determine its cost without executing the considered MRs.
- 2) For each initial input, we need to determine its coverage. In the context of metamorphic testing for Web systems, we consider input blocks based on system outputs and input parameters.
- 3) Amongst all potential input sets $I \subseteq I_{init}$, we search for a non-dominated solution I_{final} that preserves coverage while minimizing cost.

The *Automated Input Minimizer (AIM)* approach relies on analyzing the output and cost corresponding to each input. AIM obtains such information through a new feature added to the MST-wi toolset to execute each input on the system and retrieve the content of the corresponding Web pages. Obtaining the outputs of the system is very inexpensive compared to executing the considered MRs. Moreover, to address our first sub-problem, we also updated MST-wi to retrieve the cost of an input without executing the considered MRs. We rely on a surrogate metric (§ II-B), linearly correlated with execution time, which is inexpensive to collect (§ V-A).

In Step 1 (Pre-processing), AIM pre-processes the initial input set and the output information, by extracting relevant textual content from each returned Web page (§ V-B).

To address the second sub-problem, AIM relies on a *double-clustering* approach (Section VI), which is implemented by Step 2 (Output Clustering) and Step 3 (Action Clustering). For both steps, AIM relies on state-of-the-art clustering algorithms, which require to select hyper-parameter values (§ VI-A).

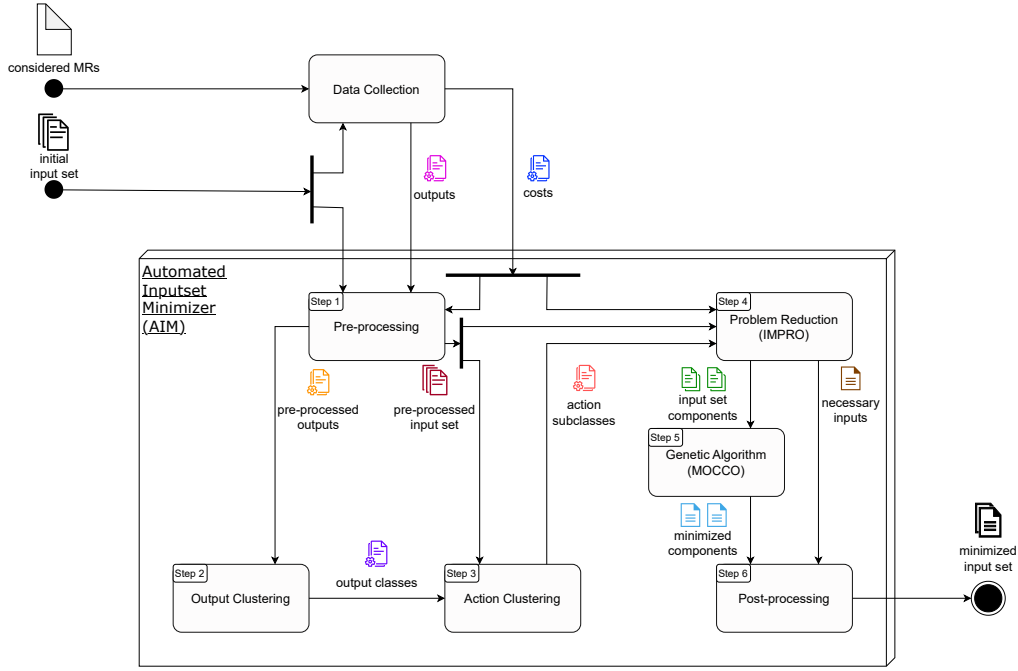


Fig. 3. Activity diagram of the Automated Input Minimizer (AIM) approach.

Output clustering (§ VI-B) is performed on the pre-processed outputs, each generated cluster corresponding to an *output class*. Then, for each output class identified by the Output Clustering Step, *Action clustering* (§ VI-C) first determines the actions whose output belongs to the considered output class, then partitions these actions based on action parameters such as URL, username, and password, obtaining *action subclasses*. On the completion of Step 3, AIM has mapped each input to a set of action subclasses, used for measuring input coverage as per our problem definition (§ III-A).

To preserve diversity in our input set, and especially to retain inputs that are necessary to exercise vulnerabilities, we require the minimized input set generated by AIM to cover the same action subclasses as the initial input set. That way, we increase the chances that the minimized input set contains at least one input able to exercise each vulnerability detectable with the initial input set.

Using cost and coverage information, AIM can address the last sub-problem. Since the size of the search space exponentially grows with the number of initial inputs, the solution cannot be obtained by exhaustive search. Actually, our problem is analogous to the knapsack problem [33], which is NP-hard, and is thus unlikely to be solved by deterministic algorithms. Therefore, AIM relies on meta-heuristic search to find a solution (Step 5) after reducing the search space (Step 4).

In Step 4, since the search space might be large, AIM first reduces the search space to the maximal extent possible (Section VII) before resorting to meta-heuristic search. Precisely, it relies on the Input set Minimization Problem Reduction Operator (IMPRO) component for problem reduction, which determines the necessary inputs, removes inputs that cannot

be part of the solution, and partition the remaining inputs into input set *components* that can be independently minimized.

In Step 5, AIM applies a genetic algorithm (Section VIII) to minimize each component. Because existing algorithms did not entirely fit our needs, we explain why we introduce MOCCO (Many-Objective Coverage and Cost Optimizer), a novel genetic algorithm which converges towards a solution covering the objectives at a minimal cost, obtaining a *minimized input set component*.

Finally, after the genetic search is completed for each component, in Step 6 (Post-processing), AIM generates the *minimized input set* by combining the necessary inputs identified by IMPRO and the inputs from the MOCCO minimized components (Section IX).

Note that, even though in our study we focus on Web systems, steps 4 (IMPRO), 5 (MOCCO), and 6 (post-processing), which form the core of our solution, are generic and can be applied to any system. Moreover, step 1 (pre-processing) and steps 2 and 3 (double-clustering) can be tailored to apply AIM to other domains (e.g., desktop applications, embedded systems). And though we relied on MST-wi to collect our data, AIM does not depend on a particular data collector, and using or implementing another data collector would enable the use of our approach in other contexts.

V. STEP 1: DATA COLLECTION AND PRE-PROCESSING

In Step 1, AIM determines the cost of each initial input (§ V-A) and extract meaningful textual content from the Web pages obtained with the initial inputs (§ V-B).

A. Input Cost

The cost of a source input is the number of actions executed by source and follow-up inputs for the considered MRs

(§ III-A). Note that counting the number of actions to be executed is inexpensive compared to executing them on the SUT, then checking for the verdict of the output relation. For instance, counting the number of actions for eleven MRs with Jenkins' initial input set (Section X) took less than five minutes, while executing the MRs took days.

B. Output Representation

Since, in this study, we focus on Web systems, the outputs of the SUT are Web pages. Fortunately, collecting these pages using a crawler and extracting their textual content is inexpensive compared to executing MRs. Hence, we can use system outputs to determine relevant input blocks (Section III).

We focus on textual content extracted from the Web pages returned by the Web system under test. We remove from the textual content of each Web page all the data that is shared among many Web pages and thus cannot characterize a specific page, like system version, date, or (when present) the menu of the Web page. Moreover, to focus on the meaning of the Web page, we consider the remaining textual content not as a string of characters but as a sequence of words. Also, following standard practice in natural language processing, we apply a stemming algorithm to consider distinct words with the same stem as equivalent, for instance the singular and plural forms of the same word. Finally, we remove stopwords, numbers, and special characters, in order to focus on essential textual information.

VI. STEPS 2 AND 3: DOUBLE CLUSTERING

To reduce the cost of MST, we want to minimize an initial input set while preserving, for each vulnerability affecting the SUT, at least one input able to exercise it; of course, in practice, such vulnerabilities are not known in advance but should be discovered by MST. Hence, we have to determine in which cases two inputs are distinct enough so that both should be kept in the minimized input set, and in which cases some inputs are redundant with the ones we already selected and thus can be removed. To determine which inputs are similar and which significantly differ, we rely on clustering algorithms. Precisely, we rely on the K-means, DBSCAN, and HDBSCAN algorithms to cluster our data points. Each of them has a set of hyper-parameters to be set and we first detail how these hyper-parameters are obtained using *Silhouette analysis* (§ VI-A).

Since, for practical reasons, we want to avoid making assumptions regarding the nature of the Web system under test (e.g., programming language or underlying middleware), we propose a black-box approach relying on input and output information to determine which inputs we have to keep or remove. In the context of a Web system, each input is a sequence of actions, each action enables a user to access a Web page (using a POST or GET request method), and each output is a Web page. After gathering output and action information, we perform *double-clustering* on our data points, i.e., two clustering steps performed in sequence:

- 1) *Output clustering* (§ VI-B) uses the outputs of the Web system under test, i.e., textual data obtained by pre-processing content from Web pages (§ V-B). We define

an output distance (§ VI-B1) to quantify similarity between these outputs, which is then used to run Silhouette analysis and clustering algorithms to partition outputs into *output classes* (§ VI-B2).

- 2) *Action clustering* (§ VI-C) then determines input coverage. First, AIM collects in the same *action set* actions leading to outputs in the same output class (§ VI-C1). Then, AIM refines each action set by partitioning the actions it contains using action parameters. To do so, it first uses the request method (§ VI-C2) to split action sets into parts. Then, we define an action distance (§ VI-C3) based on the URL (§ VI-C4) and other parameters (§ VI-C5) of the considered actions. Finally, AIM relies on Silhouette analysis and clustering algorithms to partition each part of an action set into *action subclasses* (§ VI-C6), defining our input blocks (§ III-A).

Note that *double-clustering* should not be confused with *biclustering* [34], [35], since the latter simultaneously clusters two distinct aspects (features and samples) of the data, while the former clusters only one aspect (actions, in our case, that can be seen as features) but in two consecutive steps (action outputs, then action parameters), the second refining the first one.

A. Hyper-parameters Selection

In this study, we rely on the common K-means [18], DBSCAN [19], and HDBSCAN [20] clustering algorithms (§ II-D) to determine output classes (§ VI-B) and action subclasses (§ VI-C). These clustering algorithms require a few hyper-parameters to be set. One needs to select for K-means the number of clusters k , for DBSCAN the distance threshold ϵ and the minimum number of neighbors n , and for HDBSCAN the minimum number n of individuals required to form a cluster.

To select the best values for these hyper-parameters, we rely on *Silhouette analysis*. Though the Silhouette score is a common metric used to determine optimal values for hyper-parameters [36], [37], it is obtained from the average Silhouette score of the considered data points. Thus, for instance, clusters with all data points having a medium Silhouette score cannot be distinguished from clusters where some data points have a very large Silhouette score while others have a very small one. Hence, having a large Silhouette score does not guarantee that all the data points are well-matched to their cluster. To quantify the variability in the distribution of Silhouette scores, we use Gini index, a common measure of statistical dispersion. If the Gini index is close to 0, then Silhouette scores are almost equal. Conversely, if it is close to 1, then the variability in Silhouette score across data points is large.

Hence, for our Silhouette analysis, we consider two objectives: (average) Silhouette score and the Gini index of the Silhouette scores. The selection of hyper-parameters is therefore a multi-objective problem with two objectives. We rely on the common NSGA-II evolutionary algorithm [24] to solve this problem and approximate the Pareto front regarding both Silhouette score and Gini index. Then, we select the item in the Pareto front that has the highest Silhouette score.

B. Step 2: Output Clustering

Output clustering consists in defining an output distance (§ VI-B1) to quantify dissimilarities between Web system outputs, and then to partition the outputs to obtain *output classes* (§ VI-B2).

A user communicates with a Web system using actions. Hence, an input for a Web system is a sequence of actions (e.g., login, access to a Web page, logout). As the same action may occur several times in an input in , a given occurrence of an action is identified by its position i in the input and denoted $action(in, i)$. Outputs of a Web system are textual data obtained by pre-processing the content from Web pages (§ V-B). The accessed Web page depends not only on the considered action, but also on the previous ones; for instance if the user has logged into the system. Hence, we denote by $output(in, i)$ the output of the action at position i in in .

1) *Output distance*: In this study, we use system outputs (i.e., Web pages) to characterize system states. Hence, two actions that do not lead to the same output should be considered distinct because they bring the system into different states. More generally, dissimilarity between outputs is quantified using an *output distance*. Since we deal with textual data, we consider both Levenshtein and bag distances. Levenshtein distance is usually a good representation of the difference between two textual contents [38], [39]. However, computing the minimal number of edits between two strings can be costly, since the complexity of the Levenshtein distance between two strings is $O(len(s_1) \times len(s_2))$, where $len(\cdot)$ is the length of the string [40]. Thus, we consider the bag distance [41] as an alternative to the Levenshtein distance, because its complexity is only $O(len(s_1) + len(s_2))$ [42]. But it does not take into account the order of words and is thus less precise than Levenshtein distance.

2) *Output Classes*: We partition the textual content we obtained from Web pages (§ V-B) using the K-means, DBSCAN, and HDBSCAN clustering algorithms, setting the hyperparameters using Silhouette analysis (§ VI-A), and determining similarities between outputs using the chosen output distance. We call *output classes* the obtained clusters and we denote by $OutputClass(in, i)$ the unique output class $output(in, i)$ belongs to.

C. Step 3: Action Clustering

Exercising all the Web pages is not sufficient to discover all the vulnerabilities; indeed, vulnerabilities might be detected through specific combinations of parameter values associated to an action (e.g., values belonging to a submitted form). Precisely, actions on a Web system can differ with respect to a number of *parameters* that include the URL (allowing the action to perform a request to a Web server), the method of sending a request to the server (like GET or POST), URL parameters (e.g., `http://myDomain.com/myPage?urlParameter1=value1&urlParameter2=value2`), and entries in form inputs (i.e., textarea, textbox, options in select items, datalists).

Based on the obtained output classes, *action clustering* first determines *action sets* (§ VI-C1). Then, action clustering refines each action set by partitioning the actions it contains using actions parameters. First, we give priority to the method used to send a request to the server, so we split each action set using the request method (§ VI-C2). Then, to quantify the dissimilarity between two actions, we define an action distance (§ VI-C3) based on URL (§ VI-C4) and other parameters (§ VI-C5). That way, action clustering refines each action set into *action subclasses* (§ VI-C6).

1) *Action Sets*: Based on the obtained output classes (§ VI-B2), AIM determines *action sets* such that actions leading to outputs in the same output class $outCl$ are in the same action set:

$$ActionSet(outCl) \stackrel{\text{def}}{=} \{act \mid \exists in, i : action(in, i) = act \wedge OutputClass(in, i) = outCl\}$$

Note that, because an action can have different outputs depending on the considered input, it is possible for an action to belong to several actions sets, corresponding to several output classes.

2) *Request Partition*: Each action uses either a POST or GET method to send an HTTP request to the server. Actions (such as login) that send the parameters to the server in the message body use the POST method, while actions that send the parameters through the URL use the GET method. As this difference is meaningful for distinguishing different action types, we split each action set into two parts: the actions using a POST method and those using a GET method.

3) *Action Distance*: After request partition (§ VI-C2), we consider one part of an action set at a time and we refine it using an action distance quantifying dissimilarity between actions based on remaining parameters (e.g., URL or form entries). In the context of a Web system, each Web content is identified by its URL, so we give more importance to this parameter. We denote $url(act_i)$ the URL of action act_i . For the sake of clarity, we call in the rest of the section *residual parameters* the parameters of an action which are not its request method nor its URL and we denote $res(act_i)$ the residual parameters of action act_i . Since we give more importance to the URL, we represent the distance between two actions by a real value, where the integral part corresponds to the distance between their respective URLs and the decimal part to the distance between their respective residual parameters:

$$actionDist(act_1, act_2) \stackrel{\text{def}}{=} urlDist(url(act_1), url(act_2)) + paramDist(res(act_1), res(act_2))$$

where the URL distance $urlDist(\cdot, \cdot)$ is defined in § VI-C4 and returns an integer, and the parameter distance $paramDist(\cdot, \cdot)$ is defined in § VI-C5 and returns a real number between 0 and 1.

4) *URL distance*: A URL is represented as a sequence of at least two words, separated by `://` between the first and second word, then by `/` between any other words. The length of a URL url is its number of words, denoted $len(url)$. Given two URLs, url_1 and url_2 , their *lowest common ancestor* is the longest prefix they have in common, denoted $LCA(url_1, url_2)$.

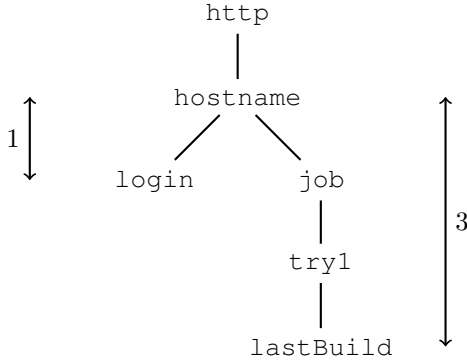


Fig. 4. The URL distance between `http://hostname/login` and `http://hostname/job/try1/lastBuild` is $1 + 3 = 4$.

We define the distance between two URLs as the total number of words separating them from their lowest common ancestor:

$$\text{urlDist}(\text{url}_1, \text{url}_2) \stackrel{\text{def}}{=} \text{len}(\text{url}_1) + \text{len}(\text{url}_2) - 2 \times \text{len}(\text{LCA}(\text{url}_1, \text{url}_2))$$

We provide an example in Figure 4.

5) *Parameter Distance*: To quantify the dissimilarity between residual parameters, we first independently quantify the dissimilarity between pairs of parameters of the same type. Since, in our context, we exercise vulnerabilities by only using string or numerical values, we ignore parameter values of other types such as byte arrays (e.g., an image uploaded to the SUT). In other contexts, new parameter distance functions suited to other input types may be required. For strings, we use the Levenshtein distance [38], [39], whereas for numerical values we consider the absolute value of their difference [43]:

$$\text{paramValDist}(v_1, v_2) \stackrel{\text{def}}{=} \begin{cases} \text{LevenshteinDist}(v_1, v_2) & \text{if } \text{type}(v_1) = \text{str} = \text{type}(v_2) \\ |v_1 - v_2| & \text{if } \text{type}(v_1) = \text{int} = \text{type}(v_2) \\ \text{undefined} & \text{otherwise} \end{cases}$$

Since we have parameters of different types, we normalize the parameter distance using the normalization function $\omega(x) = \frac{x}{x+1}$ (§ III-C). Then, we add these normalized distances together, and normalize the sum to obtain a result between 0 and 1. We compute the parameter distance in case of *matching parameters*, i.e., the number of parameters is the same and the corresponding parameters have the same type. Otherwise, we assume the largest distance possible, which is 1 due to the normalization. This is the only case where the value 1 is reached, as distance lies otherwise in $[0, 1[$, as expected for a decimal part (§ VI-C3):

$$\text{paramDist}(\text{resids}_1, \text{resids}_2) \stackrel{\text{def}}{=} \begin{cases} \omega\left(\sum_{0 \leq i < \text{len}(\text{resids}_1)} \omega(\text{paramValDist}(\text{resids}_1^{[i]}, \text{resids}_2^{[i]}))\right) & \text{if } \text{resids}_1 \text{ and } \text{resids}_2 \text{ have matching parameters} \\ 1 & \text{otherwise} \end{cases}$$

where $\text{resids}_1 = \text{res}(\text{act}_1)$, $\text{resids}_2 = \text{res}(\text{act}_2)$, and $\text{resids}^{[i]}$ is the i -th element of resids .

For instance, we consider two actions act_1 and act_2 having matching parameters with the values in Table I for page number, username, and password. The dis-

TABLE I
VALUES FOR THE EXAMPLE OF PARAMETER DISTANCE

	Page Number	Username	Password
resids_1	10	"John"	"qwerty"
resids_2	42	"Johnny"	"qwertyuiop"

tance for the page number is $\text{paramValDist}(10, 42) = 32$, normalized into $\frac{32}{32+1} \approx 0.97$. For the username, it is $\text{paramValDist}(\text{"John"}, \text{"Johnny"}) = 2$, normalized into $\frac{2}{2+1} \approx 0.66$. For the password, it is $\text{paramValDist}(\text{"qwerty"}, \text{"qwertyuiop"}) = 4$, normalized into $\frac{4}{4+1} = 0.80$. Thus, the parameter distance is $\text{paramDist}(\text{resids}_1, \text{resids}_2) \approx \frac{0.97+0.66+0.80}{0.97+0.66+0.80+1} \approx 0.71$.

6) *Action Subclasses*: We partition both parts of each action set (§ VI-C2) using the K-means, DBSCAN, or HDBSCAN clustering algorithms, setting the hyper-parameters using our Silhouette analysis (§ VI-A), and quantifying action dissimilarity using our action distance (§ VI-C3), obtaining clusters we call *action subclasses*. We denote by $\text{ActionSubclass}(\text{act}, \text{actSet})$ the unique action subclass bl from the action set actSet such that $\text{act} \in bl$. For the sake of simplicity, we denote $\text{Subclass}(\text{in}, i)$ the action subclass corresponding to the i -th action in input in :

$$\text{Subclass}(\text{in}, i) \stackrel{\text{def}}{=} \text{ActionSubclass}(\text{action}(\text{in}, i), \text{ActionSet}(\text{OutputClass}(\text{in}, i)))$$

Finally, in our study, the objectives covered by an input are:

$$\text{Cover}(\text{in}) \stackrel{\text{def}}{=} \{\text{Subclass}(\text{in}, i) \mid 1 \leq i \leq \text{len}(\text{in})\}$$

VII. STEP 4: PROBLEM REDUCTION

The search space for our problem (§ III-D) consists of all the subsets of the initial input set, which leads to 2^m potential solutions, where m is the number of initial inputs.

For this reason, AIM integrates a *problem reduction* step, implemented by the Input set Minimization Problem Reduction Operator (IMPRO) component, to minimize the search space before solving the search problem in the next step (Section VIII). We apply the following techniques to reduce the size of the search space:

- *Determining redundancy*: Necessary inputs (§ III-C) must be part of the solution, hence one can only investigate redundant inputs (§ VII-C). Moreover, one can restrict the search by removing the objectives already covered by necessary inputs. Finally, if a redundant input does not cover any of the remaining objectives, it will not contribute to the final coverage, and hence can be removed.
- *Removing duplicates*: Several inputs may have the same cost and coverage. In this case, we consider them as *duplicates* (§ VII-D). Thus, we keep only one and we remove the others.
- *Removing locally-dominated inputs*: For each input, if there exists other inputs that cover the same objectives at a same or lower cost, then the considered input is *locally-dominated* by the other inputs (§ VII-E) and is removed.
- *Dividing the problem*: We consider two inputs covering a common objective as being connected. Using this

relation, we partition the search space into connected *components* that can be independently solved (§ VII-F), thus reducing the number of objectives and inputs to investigate at a time.

Before detailing these techniques, we first explain in which order there are applied (§ VII-A).

A. Order of Reduction Techniques

We want to perform first the least expensive reduction techniques, to sequentially reduce the cost of the following more expensive techniques. Determining redundancy requires $\mathcal{O}(m \times c)$ steps, removing duplicates requires $\mathcal{O}(m^2)$ steps, and removing locally-dominated inputs requires $\mathcal{O}(m \times 2^n)$ steps, where m is the number of inputs, c is the maximal number of objectives covered by an input, and n is the maximal number of neighbors for an input (i.e., the number of other inputs that cover an objective shared with the considered input). In our study, we assume $c < m$. Hence, we first determine redundancy, then remove duplicates, and remove locally-dominated inputs. Dividing the problem requires exploring neighbors and comparing non-visited inputs with visited ones, so it is potentially the most costly of these reduction techniques; hence, it is performed at the end.

After determining redundancy, the removal of already covered objectives may lead to new inputs being duplicates or locally-dominated. Moreover, the removal of duplicates or locally-dominated inputs may lead to changes in redundancy, making some previously redundant inputs necessary. Hence, these reduction techniques should be iteratively applied, until a stable output is reached. Such output can be detected by checking if inputs were removed during an iteration.

Therefore, the order is as follows. We first initialize variables (§ VII-B). Then we repeat, until no input is removed, the following steps: determine redundancy (§ VII-C), remove duplicates (§ VII-D), and remove locally-dominated inputs (§ VII-E). Finally, we divide the problem into sub-problems (§ VII-F).

B. Initializing Variables

During problem reduction, we consider three variables: I_{necess} , the set of the inputs that has to be part of the final solution, I_{search} , the remaining inputs to be investigated, and $Coverage_{obj}$, the objectives that remain to be covered by subsets of I_{search} . I_{necess} is initially empty. I_{search} is initialized as the initial input set. $Coverage_{obj}$ is initialized as the coverage of the initial input set (§ VI-C6 and § III-A).

C. Determining Redundancy

This technique is presented in Algorithm 1. Each time it is repeated, the redundancy of the remaining inputs is computed (Line 2). Among them, inputs which are necessary (§ III-C) in I_{search} (Line 3) for the objectives in $Coverage_{obj}$ have to be included in the final solution (Line 4), otherwise some objectives will not be covered. Then, the objectives already covered by the necessary inputs are removed (Line 5). Hence, in the following, we only consider, for each remaining input

Algorithm 1 Redundancy determination technique.

```

1: procedure REDUNDANCY( $I_{necess}, I_{search}, Coverage_{obj}$ )
2:    $I_{redund} \leftarrow Redundant(I_{search})$ 
3:    $I_{necess}^{new} \leftarrow I_{search} \setminus I_{redund}$ 
4:    $I_{necess} \leftarrow I_{necess} \cup I_{necess}^{new}$ 
5:    $Coverage_{obj} \leftarrow Coverage_{obj} \setminus Cover(I_{necess}^{new})$ 
6:    $I_{search} \leftarrow \{in \in I_{redund} \mid Cover(in) \cap Coverage_{obj} \neq \emptyset\}$ 
7:   return  $I_{necess}, I_{search}, Coverage_{obj}$ 

```

$in \in I_{search}$, their coverage regarding the remaining objectives, i.e., $Cover(in) \cap Coverage_{obj}$, instead of $Cover(in)$.

Finally, some redundant inputs may cover only objectives that are already covered by necessary inputs. In that case, they cannot be part of the final solution because they would contribute to the cost but not to the coverage of the objectives. Hence, we restrict without loss the search space for our problem by considering only redundant inputs that can cover the remaining objectives (Line 6).

D. Removing Duplicates

In the many-objective problem described in § III-C, inputs are characterized by their coverage (§ VI-C6) and their cost (§ III-A). Hence, two inputs with the same coverage and cost are considered duplicates. In that case, IMPRO selects one and remove the other.

E. Removing Locally-dominated Inputs

For a given input $in \in I_{search}$, if the same coverage can be achieved by one or several other input(s) for at most the same cost, then in is not required for the solution. Formally, we say that the input $in \in I_{search}$ is *locally dominated* by the subset $S \subseteq I_{search}$, denoted $in \sqsubseteq S$, if $in \notin S$, $Cover(in) \subseteq Cover(S)$, and $cost(in) \geq cost(S)$. In order to simplify the problem, inputs that are locally dominated should be removed from the remaining inputs I_{search} .

Removing a redundant input in (§ III-C) can only affect the redundancy of the inputs in I that cover objectives in $Cover(in)$. Hence, we consider two inputs as being connected if they cover at least one common objective. Formally, we say that two inputs in_1 and in_2 *overlap*, denoted by $in_1 \sqcap in_2$, if $Cover(in_1) \cap Cover(in_2) \neq \emptyset$. The name of the *local* dominance relation comes from the fact proved in the separate appendix (Proposition 3) that, to determine if an input is locally-dominated, one has only to check amongst its neighbors for the overlapping relation instead of amongst all the remaining inputs, thus making this step tractable.

One concern is that removing a locally-dominated input could alter the local dominance of other inputs. Fortunately, this is not the case for local dominance. We prove in the separate appendix (Theorem 2) that, for every locally-dominated input $in \in I_{search}$, there always exists a subset $S \subseteq I_{search}$ of not locally-dominated inputs such that $in \sqsubseteq S$. Hence, the locally dominated inputs can be removed in any order without reducing coverage or preventing cost reduction, both being ensured by non locally-dominated inputs. Therefore, IMPRO

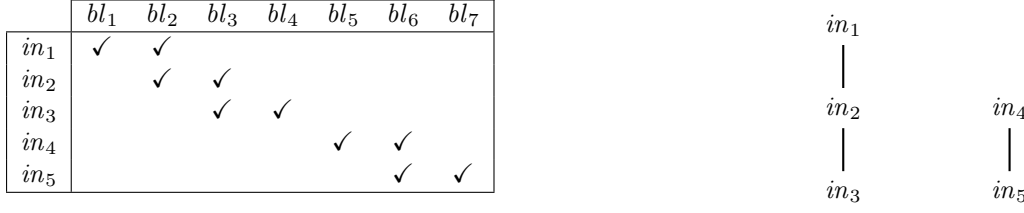


Fig. 5. Inputs covering one objective in common (left) are connected in the corresponding overlapping graph (right).

keeps in the search space only the remaining inputs that are not locally-dominated:

$$I_{search} \leftarrow \{in \in I_{search} \mid \forall S \subseteq I_{search} : in \not\subseteq S\}$$

F. Dividing the Problem

After removing as many inputs as possible, we leverage the overlapping relation (§ VII-E) to partition the remaining inputs into connected components, in a divide-and-conquer approach. We denote $\mathcal{G}_{\square}(I)$ the *overlapping graph* of the input set I , i.e., the undirected graph such that vertices are inputs in I and edges correspond to the overlapping relation \square , and $Comps(I)$ the set of the *connected components* of $\mathcal{G}_{\square}(I)$.

For instance, in Figure 5, we represent on the left the input blocks covered by each input in the search space and the corresponding overlapping graph on the right. The connected components of the graph are $\{in_1, in_2, in_3\}$ and $\{in_4, in_5\}$.

Such connected components are important because we prove in the separate appendix (Proposition 8) that inputs in a connected component can be removed without altering the redundancy of inputs in other connected components. Similarly, we prove in the separate appendix (Theorem 3) that the gain, i.e., the maximal cost reduction from removing redundant inputs (§ III-C), can be independently computed on each connected component, i.e., for each input set I , $gain(I) = \sum_{C \in Comps(I)} gain(C)$.

Hence, instead of searching a solution on I_{search} to solve our initial problem (§ III-D), we use a divide-and-conquer strategy to split the problem into more manageable sub-problems that can be independently solved on each connected component $C \in Comps(I_{search})$.

We denote $Coverage_{obj}(C) \stackrel{\text{def}}{=} Coverage_{obj} \cap Cover(C)$ the remaining objectives to be covered by inputs in C and we formulate the sub-problem on the connected component similarly to the initial problem:

$$\text{minimize}_{I \subseteq C} F_C(I) \stackrel{\text{def}}{=} [\omega(cost(I)), f_{bl_1}(I), \dots, f_{bl_n}(I)]$$

where $Coverage_{obj}(C) = \{bl_1, \dots, bl_n\}$ and the minimize notation is detailed in § II-E1. We denote I_C a non-dominated solution with full coverage, such that $F_C(I_C) = [\omega(cost_{\min}), 0, \dots, 0]$.

VIII. STEP 5: GENETIC SEARCH

Our initial goal of obtaining a subset $I_{final} \subseteq I_{init}$ with total input coverage at minimal cost (§ III-A) can be expressed as a weighted set cover problem. Given a universe \mathcal{U} , a

set $S = \{S_1, S_2, \dots\}$ of subsets $S_i \subseteq \mathcal{U}$, and a weight function mapping each subset S_i to a positive real number, the weighted set cover problem consists in finding a subset $T \subseteq S$ such that subsets S_i in T cover \mathcal{U} at a minimal cumulative cost. To express our initial goal as a weighted set cover problem, we consider as universe the input blocks $\mathcal{U} = Cover(I_{init})$ to cover, as subsets of the universe the input blocks $Cover(in)$ covered by each input $in \in I_{init}$, and as weight of a subset $Cover(in)$ the cost $cost(in)$ of the corresponding input. A solution to this set cover problem is a set of subsets $T = \{Cover(in_1), \dots, Cover(in_n)\}$, providing a minimized input set $\{in_1, \dots, in_n\}$, given that each subset $Cover(in)$ can be uniquely mapped to its initial input in . This is the case for each connected component C , since IMPRO removed duplicates (§ VII-D) and locally-dominated inputs (§ VII-E). While the set cover problem is NP-hard [44], it admits a very efficient [45] polynomial time approximation using a greedy algorithm, as well as various techniques to find approximate solutions to an ILP formulation of the problem [46]. Unfortunately, the set cover problem considers only the total coverage and the cumulative cost of solutions, not input blocks as individual objectives (§ III-B) as in the formulation of our many-objective problem (§ III-D) or sub-problem (§ VII-F). Hence, an algorithm solving the set cover problem may not take into account relevant information about how each input block is covered. Indeed, some blocks may be more difficult to cover and should thus receive priority. Further, an input selected by the greedy algorithm, because it covers many blocks at small cost, may be less optimal than several inputs covering more blocks at slightly larger cost. For instance, if $Cover(in_1) = \{bl_1, bl_2\}$ with $cost(in_1) = 2$, $Cover(in_2) = \{bl_1, bl_3\}$ with $cost(in_2) = 3$, and $Cover(in_3) = \{bl_2, bl_4\}$ with $cost(in_3) = 3$, the greedy algorithm selects in_1 because it is locally optimal, then in_2 and in_3 to cover respectively bl_3 and bl_4 , at a total cost of 8. However, selecting only in_2 and in_3 leads to full coverage but at a lower cost of 6. In this example, the greedy algorithm selects in_1 covering blocks bl_1 and bl_2 , which are easier to cover because each block can be covered by two inputs, while blocks bl_3 and bl_4 can only be covered by one input. Therefore, it is unlikely for a polynomial time approximation of the set cover problem to find, in general, the best solution to our many-objective problem.

Alternatively, obtaining an optimal solution I_C to our sub-problem on a connected component C is similar to solving the knapsack problem, which is also NP-hard [33]. To be more precise, our problem is equivalent to the 0-1 knapsack prob-

lem, which consists in selecting a subset of items to maximize a total value, while satisfying a weight capacity. Since we consider a many-objective problem (§ III-B), we must address the multidimensional variant of the 0-1 knapsack problem, where each item has many “weights”, one per considered objective. In our case, we minimize the total cost instead of maximizing total value and ensure the coverage of each action objective instead of making sure that each weight capacity is not exceeded. Furthermore, the 0-1 multidimensional knapsack problem is harder than the initial knapsack problem as it does not admit a fully polynomial time approximation scheme [47], hence the need for a meta-heuristic.

For our approach to scale, we adopt a genetic algorithm because it is known to find good approximations in reasonable execution time [28] and has been widely used in software testing. An input set $I \subseteq C$ can thus be seen as a chromosome, where each gene corresponds to an input $in \in C$, the gene value being 1 if $in \in I$ and 0 otherwise. Though several many-objective algorithms have been successfully applied within the software engineering community, like NSGA-III [27], [29] and MOSA [28] (§ II-E2), these algorithms do not entirely fit our needs (§ VIII-A). Hence, we propose MOCCO (Many-Objective Coverage and Cost Optimizer), a novel genetic algorithm based on two populations and summarized in Algorithm 2. We first explain how these populations are initialized (§ VIII-B). Then, for each generation, MOCCO performs the following standard steps: selection of the parents (§ VIII-C), crossover of the parents to produce an offspring (§ VIII-D), mutation of the offspring (§ VIII-E), and update of the populations (§ VIII-F) to obtain the next generation. The process continues until a termination criterion is met (§ VIII-G). Then, we detail how MOCCO determines the solution I_C to each sub-problem.

A. Motivation for a Novel Genetic Algorithm

While our problem is similar to the multidimensional 0-1 knapsack problem, it is not exactly equivalent, since standard solutions to the multidimensional knapsack problem have to ensure that, for each “weight type”, the total weight of the items in the knapsack is below weight capacity, while we want to ensure that, for each objective, at least one input in the minimized input set covers it. Hence, standard solutions based on genetic search are not applicable in our case, and we focus on genetic algorithms able to solve many-objective problems in the context of test case generation or minimization. We have explained earlier (§ II-E2) the challenges raised by many-objective problems and how the NSGA-III [27], [29] and MOSA [28] genetic algorithms tackle such challenges.

NSGA-III has the advantage of letting users choose the parts of the Pareto front they are interested in, by providing reference points. Otherwise, it relies on a systematic approach to place points on a normalized hyperplane. While this approach is useful in general, we are interested only in solutions that are close to a utopia point $[0, 0, \dots, 0]$ (§ III-C) covering every objective at no cost. Hence, we do not care about diversity over the Pareto front, and we want to explore a very specific region of the search space. Moreover, apart from the starting

Algorithm 2 MOCCO overview.

```

1: procedure MOCCO( $C, n_{size}, n_{gens}, time_{budget}$ )
2:    $time_{start} \leftarrow getTime()$ 
3:    $Roofers \leftarrow initRoofers(C, n_{size})$ 
4:    $Misers \leftarrow \emptyset$ 
5:    $n \leftarrow 1$ 
6:    $stillTime \leftarrow True$ 
7:   while  $n \leq n_{gens} \wedge stillTime$  do
8:      $n \leftarrow n + 1$ 
9:      $I_1, I_2 \leftarrow selectParents(Roofers, Misers)$ 
10:     $I_3, I_4 \leftarrow crossover(I_1, I_2)$ 
11:    for  $I \in \{I_3, I_4\}$  do
12:       $I \leftarrow mutate(I)$ 
13:       $I \leftarrow reduce(I)$ 
14:      if  $getTime() - time_{start} > time_{budget}$  then
15:         $stillTime \leftarrow False$ 
16:      break
17:      if  $I \in Roofers \cup Misers$  then
18:        continue
19:      if  $Cover(I) = Coverage_{obj}(C)$  then
20:         $Roofers \leftarrow updRoofers(Roofers, I)$ 
21:      else
22:         $Misers \leftarrow updMisers(Misers, I)$ 
23:     $I_C \leftarrow selectSolution(Roofers)$ 
24:  return  $I_C$ 

```

points, the main use of the reference points in NSGA-III is to determine, after the first nondomination fronts are obtained from NSGA-II [24], the individuals to be selected from the last considered front, so that the population reaches a predefined number of individuals. This is not a problem we face because we know there is only one point (or several points, but at the same coordinates) in the Pareto front that would satisfy our constraint of full coverage at minimal cost. Hence, we do not use the Pareto front as a set of solutions, even if we intend to use individuals in the Pareto front as intermediate steps to reach relevant solutions.

Regarding MOSA [28], trade-offs obtained from approximating the Pareto front are only used for maintaining diversity during the search, which is similar to what we intend to do. But, as opposed to the use case tackled by MOSA, in our case determining inputs covering a given objective is straightforward. Indeed, for each objective, we can easily determine inputs that are able to cover it (§ III-B). Hence, individuals ensuring the coverage of the objectives are easy to obtain, while the hard part of our problem is to determine a combination of inputs able to cover all the objectives at a minimal cost. Hence, even if MOSA may find a reasonable solution, because it focuses on inputs individually covering an objective and not on their collective coverage and cost, it is unlikely to find the best solution.

Hence, we propose a novel genetic algorithm, named MOCCO. We take inspiration from MOSA [28] by considering two populations: 1) a population of solutions (like MOSA’s archive), called the *roofers* because they cover all the objectives (§ VI-C6), and 2) a population of individuals on the

Pareto front (§ III-D), called the *misers* because they minimize the cost, while not covering all objectives. Like NSGA-III and MOEA, MOCCO has to tackle challenges raised by many-objective problems (§ II-E2).

To address challenge 1, we take inspiration from the whole suite approach [31] which counts covered branches as a single objective, by defining the *exposure* as the sum of the coverage objective functions (§ III-C):

$$exposure(I) \stackrel{\text{def}}{=} \sum_{bl_i \in Coverage_{obj}(C)} f_{bl_i}(I)$$

Since $f_{bl_i}(\cdot)$ is zero when the objective bl_i is covered, the larger the exposure, the smaller the input set coverage. As described in § III-B, we do not use the exposure as objective because we want to distinguish between input blocks. But we use it as a weight when randomly selecting a parent amongst the misers (§ VIII-C), so that the further away a miser is from complete coverage, the less likely it is to be selected. That way, we aim to benefit from the large number of dimensions to avoid getting stuck in a local optimum and to have a better convergence rate [23], while still focusing the search on the region of interest.

Since we want to deeply explore this particular region, we do not need to preserve diversity over the whole Pareto front. Therefore, we do not use diversity operators, avoiding challenge 2.

Finally, we address challenge 3 by 1) restricting the recombination operations and 2) tailoring them to our problem, as follows:

- 1) A crossover between roofers can only happen during the first generations, when no miser is present in the population. After the first miser is generated, crossover (§ VIII-D) is allowed only between a roofer and a miser. Hence, the roofer parent provides full coverage while the miser parent provides almost full coverage at low cost. Moreover, because of how the objective functions are computed (§ III-C), the not-yet-covered objectives are likely to be covered in an efficient way. That way, we hope to increase our chances of obtaining offspring with both large coverage and low cost.
- 2) Not only our recombination strategy is designed to be computationally efficient (by minimizing the number of introduced redundancies), but we exploit our knowledge of input coverage to determine a meaningful crossover between parents, with inputs from one parent for one half of the objectives and inputs from the other parent for the other half.

B. Population Initialization

During the search, because we need diversity to explore the search space, we consider a population (with size $n_{size} \geq 2$) of the least costly individuals generated so far that satisfy full coverage. We call *roofers* such individuals, by analogy with covering a roof, and we denote $Roofers(n)$ the roofer population at generation n .

But focusing only on the roofers would prevent us to exploit the least expensive solutions obtained in the Pareto front

Algorithm 3 Roofer population initialization.

```

1: procedure INITROOFERS( $C, n_{size}$ )
2:    $Roofers \leftarrow \emptyset$ 
3:   while  $|Roofers| < n_{size}$  do
4:      $I \leftarrow \emptyset$ 
5:     while  $Cover(I) \neq Coverage_{obj}(C)$  do
6:        $bl \leftarrow select(Coverage_{obj}(C) \setminus Cover(I), P_{unif})$ 
7:        $in \leftarrow select(Inputs(bl), P_{init})$ 
8:        $I \leftarrow I \cup \{in\}$ 
9:        $I \leftarrow reduce(I)$ 
10:    if  $I \notin Roofers$  then
11:       $Roofers \leftarrow Roofers \cup \{I\}$ 
12:  return  $Roofers$ 

```

while trying to minimize the connected component (§ VII-F). Instead, inputs that do not cover all the objectives, and will thus not be retained for the final solution, but are efficient at minimizing cost, are thus useful as intermediary steps towards finding an efficient solution. Hence, we maintain a second population, formed by individuals that are non-dominated so far and minimize cost while not covering all the objectives. We call *misers* such individuals, because they focus on cost reduction more than objective coverage, and we denote $Misers(n)$ the miser population at generation n .

The reason for maintaining two distinct populations is to restrict the crossover strategy (§ VIII-D) so that (in most cases) one parent is a roofer and one parent is a miser. Since misers prioritize cost over coverage, a crossover with a miser tends to reduce cost. Because roofers prioritize coverage over cost, a crossover with a roofer tends to increase coverage. Hence, with such a strategy, we intend to converge towards a solution minimizing cost and maximizing coverage.

For both populations, we want to ensure that the individuals are reduced (§ III-C), i.e., they contain no redundant inputs. Hence, during the initialization and updates of these populations, we ensure that removal steps are performed. Because, as detailed in the following, the number of redundant inputs obtained for each generation is small, the optimal order of removal steps can be exhaustively computed. We denote by $reduce(I)$ the input set I after these removal steps. This limits the exploration space, since non-reduced input sets are likely to have a large cost and hence to be far away for the utopia point of full coverage at no cost we intend to focus on.

The miser population is initially empty, i.e., $Misers(0) \stackrel{\text{def}}{=} \emptyset$, as misers are generated during the search through mutations (§ VIII-F). We detail in Algorithm 3 how the roofer population $Roofers(0)$ is initialized, where $select(X, P)$ randomly select one element in X using distribution P , P_{unif} denotes the uniform distribution, and P_{init} denotes the following distribution:

$$P_{init}(in_1) \stackrel{\text{def}}{=} \frac{\frac{1}{1 + occurrence(in_1)}}{\sum_{in_2 \in Inputs(bl)} \frac{1}{1 + occurrence(in_2)}}$$

where $occurrence(in)$ denotes the number of times the input in was selected in the roofer population so far. This distribution ensures that inputs that were not selected so far are more

likely to be selected, so that the initial roofer population can be more diverse. Note that computing $reduce(I)$ is tractable since adding a new input can only affect the redundancy of inputs that overlap with it (§ VII-E).

C. Parents Selection

For each generation n , parents are selected as follows. If $Misers(n) \neq \emptyset$, then one parent is selected from the miser population and one from the roofer population. Otherwise, two distinct parents are selected from the roofer population. A parent $I_1 \in Misers(n)$ is randomly selected from the miser population using the following distribution:

$$P_{misers}(I_1) \stackrel{\text{def}}{=} \frac{\frac{1}{exposure(I_1)}}{\sum_{I_2 \in Misers(n)} \frac{1}{exposure(I_2)}}$$

where the exposure is defined in § VIII-A. The purpose of this distribution is to ensure that input sets with large coverage or at least large potential (§ III-C) are more likely to be selected. A parent $I_1 \in Roofers(n)$ is randomly selected from the roofer population using the following distribution:

$$P_{roofers}(I_1) \stackrel{\text{def}}{=} \frac{\frac{1}{cost(I_1)}}{\sum_{I_2 \in Roofers(n)} \frac{1}{cost(I_2)}}$$

The purpose of this distribution is to ensure that less costly input sets are more likely to be selected.

D. Parents Crossover

After selecting two distinct parents I_1 and I_2 , we detail how they are used to generate the offspring I_3 and I_4 . Our crossover strategy exploits the fact that, for each objective bl to cover, it is easy to infer inputs in $Inputs(bl)$ able to cover bl (§ III-B). For each crossover, we randomly split the objectives in two halves O_1 and O_2 such that $O_1 \cup O_2 = Coverage_{obj}(C)$ and $O_1 \cap O_2 = \emptyset$. We consider here a balanced split, to prevent cases where one parent massively contributes to offspring coverage.

Then, we use this split to define the crossover: inputs in the connected component C are split between $S_1 \stackrel{\text{def}}{=} Inputs(O_1)$, the ones covering the first half of the objectives, and $S_2 \stackrel{\text{def}}{=} Inputs(O_2)$, the ones covering the second half. Note that some inputs may cover objectives both in O_1 and O_2 , so we call the *edge* of the split the intersection $S_1 \cap S_2$. Because we assume both parents are reduced, this means that redundant inputs can only happen at the edge of the split. The genetic material of both parents I_1 and I_2 is then split in two parts: inputs in S_1 and inputs in S_2 , as follows:

$$\begin{aligned} I_1 &= (I_1 \cap S_1) \cup (I_1 \cap S_2) \\ I_2 &= (I_2 \cap S_1) \cup (I_2 \cap S_2) \end{aligned}$$

Then, these parts are swapped to generated the offspring, as follows:

$$\begin{aligned} I_3 &\stackrel{\text{def}}{=} (I_1 \cap S_1) \cup (I_2 \cap S_2) \\ I_4 &\stackrel{\text{def}}{=} (I_2 \cap S_1) \cup (I_1 \cap S_2) \end{aligned}$$

For illustration purpose we consider in Figure 6 a small connected component $C = \{in_1, in_2, in_3, in_4, in_5\}$ and the

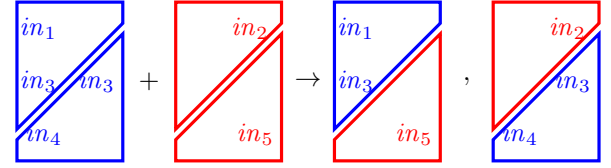


Fig. 6. Crossover Example

following split for the inputs: $Inputs(O_1) = \{in_1, in_2, in_3\}$ and $Inputs(O_2) = \{in_3, in_4, in_5\}$. The edge of the split is $Inputs(O_1) \cap Inputs(O_2) = \{in_3\}$. The offspring of parents $I_1 = \{in_1, in_3, in_4\}$ and $I_2 = \{in_2, in_5\}$ is $I_3 = \{in_1, in_3, in_5\}$ and $I_4 = \{in_2, in_3, in_4\}$.

E. Offspring Mutation

Each gene of an offspring I corresponds to an input $in \in C$, the gene value being 1 if $in \in I$ and 0 otherwise. A mutation happens when this gene value is changed, hence mutation randomly adds or removes one input from an offspring.

The crossover (at the edge of the split) and mutation (when an input is added) steps may result in inputs being redundant in the offspring. Since redundancies in the connected component were already reduced by IMPRO (§ VII-F) and changes in redundancy could happen only amongst neighbors (for the overlapping relation) of the changed inputs (§ VII-E), we only expect a few redundant inputs. Therefore, removal steps can be exhaustively computed to replace each offspring I by its reduced counterpart $reduce(I)$ (§ VIII-B).

F. Population Update

We detail how the offspring is used to obtain the populations $Roofers(n+1)$ and $Misers(n+1)$ at generation $n+1$.

First, we discard any offspring that is a duplicate of an individual already present in either $Roofers(n)$ or $Misers(n)$. Indeed, the duplication of individuals would only result in altering their weight for being selected as parents (thus, the intended procedure), reducing roofer diversity, and increasing the number of miser comparisons (as detailed below).

If a remaining offspring I_1 covers all the objectives in $Coverage_{obj}(C)$, then it is a candidate for the roofer population. Otherwise, it is a candidate for the miser population.

For each candidate I_1 for the roofer population, MOCCO computes its cost. If $cost(I_1) \leq \max\{cost(I) \mid I \in Roofers(n)\}$, then it selects the most costly roofer I_2 (or, in case of a tie, one of the most costly roofers), removes I_2 from the population, and adds I_1 . Otherwise, I_1 is rejected. Note that we chose \leq instead of $<$ for the above cost criterion because, in case of a tie, we prefer to evolve the population instead of maintaining the status quo, to increase the odds of exploring new regions of the search space.

For each candidate I_1 to the miser population, we compute its fitness vector $F_C(I_1)$ (§ VII-F). Then, for each $I_2 \in Misers(n)$ we compare $F_C(I_1)$ and $F_C(I_2)$. If $I_2 \succ I_1$ in the sense of Pareto-dominance (§ II-E1), then we stop the process and I_1 is rejected. If $I_1 \succ I_2$, then I_2 is removed from

$Misers(n)$. That way, we ensure that the miser population contains only non-dominated individuals. After completing the comparisons, if the process was not stopped, then I_1 itself is non-dominated, so it is added to the miser population. In that case, $F_C(I_1)$ is stored for future comparisons.

Properties satisfied by roofers and misers are detailed in the separate appendix (Theorems 4 and 5).

G. Termination

MOCCO repeats the process until it reaches a fixed number of generations or exhausts a given time budget. Then, amongst the least costly roofers (several may have the same cost), it randomly selects one individual I_C as solution to our sub-problem (§ VII-F). I_C covers all the objectives and, amongst the input sets covering those objectives, I_C has the smallest cost encountered during the search.

IX. STEP 6: DATA POST-PROCESSING

The set I_{necess} was initially empty (§ VII-B) and then accumulated necessary inputs each time redundancy was determined (§ VII-C). After removing inputs and reducing the objectives to be covered accordingly (Section VII), IMPRO obtained a set I_{search} of remaining inputs and objectives. Then, IMPRO divided the remaining problem into sub-problems (§ VII-F), one for each connected component C . Finally, for each connected component C , the corresponding sub-problem was solved using MOCCO (Section VIII), obtaining the corresponding minimized component I_C . At the end of the search, AIM merges inputs from each minimized component I_C with the necessary inputs I_{necess} to obtain a *minimized input set* I_{final} as solution to our initial problem (§ III-D):

$$I_{final} \stackrel{\text{def}}{=} I_{necess} \cup \bigcup_{C \in \text{Comps}(I_{search})} I_C$$

X. EMPIRICAL EVALUATION

In this section, we report our results on the assessment of our approach with two Web systems. We investigate the following Research Questions (RQs):

- RQ1 What is the vulnerability detection effectiveness of AIM, compared to alternatives?** This research question aims to determine if and to what extent AIM reduces the effectiveness of MST by comparing the vulnerabilities detected between the initial and the minimized input sets. Also, we further compare the vulnerability detection rate of AIM with simpler alternative approaches.
- RQ2 What is the input set minimization effectiveness of AIM, compared to alternatives?** This research question aims to analyze the magnitude of minimization in terms of the number of inputs, cost (§ III-A and V-A), and execution time for the considered MRs, both for AIM and alternative approaches.
- RQ3 What is the input set minimization effectiveness of MOCCO, compared to alternatives?** This research question aims to determine the particular contribution

of the novel MOCCO genetic algorithm in minimizing cost while preserving vulnerability detection, by comparing it to alternative approaches for different time budgets.

A. Experiment Design

1) Subjects of the Study: To assess our approach with MRs and input sets that successfully detect real-world vulnerabilities, we rely on the same input sets and settings as MST-wi [6].

The targeted Web systems under test are Jenkins [48] and Joomla [49]. Jenkins is a leading open source automation server while Joomla is a content management system (CMS) that relies on the MySQL RDBMS and the Apache HTTP server. We chose these Web systems because of their plug-in architecture and Web interface with advanced features (such as Javascript-based login and AJAX interfaces), which makes Jenkins and Joomla good representatives of modern Web systems.

Further, these systems present differences in their output interface and input types that, since inputs and outputs are key drivers for our approach, contribute to improve the generalizability of our results. Concerning outputs, Joomla is a CMS where Web pages tend to contain a large amount of static text that differ in every page, while Jenkins provides mainly structured content that may continuously change (e.g., seconds from the last execution of a Jenkins task). The input interfaces of Jenkins are mainly short forms and buttons whereas the inputs interfaces of Joomla often include long text areas and several selection interfaces (e.g., for tags annotation).

The selected versions of Jenkins and Joomla—2.121.1 and 3.8.7, respectively—are affected by known vulnerabilities that can be triggered from the Web interface; we describe them in § X-A2.

The input set provided in the MST-wi’s replication package has been collected by running Crawljax with, respectively, four users for Jenkins and six users for Joomla having different roles, e.g., admin. For each role, Crawljax has been executed for a maximum of 300 minutes, to prevent the crawler from running indefinitely, thereby avoiding excessive resource consumption. Further, to exercise features not reached by Crawljax, a few additional Selenium [50]-based test scripts (four for Jenkins and one for Joomla) have been added to the input set. In total, we have 160 initial inputs for Jenkins and 148 for Joomla, which are all associated to a unique identifier.

Since MOCCO assumes redundancy in the input set to be already reduced by IMPRO (e.g., to ensure the reduction step after mutation is tractable), we do not consider the initial input sets for RQ3 but their reduced versions. Thus, we run the double-clustering (Section VI) and problem reduction (Section VII) steps to obtain a set of necessary inputs and several connected components. Then, the input set for RQ3 is the union of these necessary inputs and connected components, called a *reduced input set*. We repeat this process several times (to reduce the impact of randomness), thus obtaining several reduced input sets.

Finally, since we know necessary inputs should be part of the solution and minimization can be performed independently

TABLE II
JENKINS VULNERABILITIES.

CVE	Description	Vulnerability Type	Input Identifiers
CVE-2018-1000406 [52]	In the file name parameter of a Job configuration, users with Job / Configure permissions can specify a relative path escaping the base directory. Such path can be used to upload a file on the Jenkins host, resulting in an arbitrary file write vulnerability.	CWE_22	160
CVE-2018-1000409 [53]	A session fixation vulnerability prevents Jenkins from invalidating the existing session and creating a new one when a user signed up for a new user account.	CWE_384	112, 113, 114
CVE-2018-1999003 [54]	Jenkins does not perform a permission check for URLs handling cancellation of queued builds, allowing users with Overall / Read permission to cancel queued builds.	CWE_280, CWE_863	116, 157
CVE-2018-1999004 [55]	Jenkins does not perform a permission check for the URL that initiates agent launches, allowing users with Overall / Read permission to initiate agent launches.	CWE_863, CWE_285	2, 116
CVE-2018-1999006 [56]	A exposure of sensitive information vulnerability allows attackers to determine the date and time when a plugin was last extracted.	CWE_200, CWE_668	33, 55, 57, 61, 62, 63, 64, 75, 107, 108, 110, 135, 136, 156, 160
CVE-2018-1999046 [57]	Users with Overall / Read permission are able to access the URL serving agent logs on the UI due to a lack of permission checks.	CWE_200	2, 116
CVE-2020-2162 [58]	Jenkins does not set Content-Security-Policy headers for files uploaded as file parameters to a build, resulting in a stored XSS vulnerability.	CWE_79	1, 18, 19, 23, 26, 75, 156, 158
Password aging problem in Jenkins	Jenkins does not integrate any mechanism for managing password aging; consequently, users aren't incentivized to update passwords periodically.	CWE_262	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 18, 19, 20, 22, 23, 24, 25, 26, 27, 28, 30, 32, 110, 33, 34, 35, 38, 39, 41, 42, 43, 44, 45, 46, 47, 58, 61, 62, 64, 65, 66, 69, 70, 71, 73, 74, 75, 104, 108, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 133, 134, 143, 145, 146, 159, 160
Weak password in Jenkins	Jenkins does not require users to have strong passwords, which makes it easier for attackers to compromise user accounts.	CWE_521	112, 113, 114

TABLE III
Joomla VULNERABILITIES.

CVE	Description	Vulnerability Type	Input Identifiers
CVE-2018-11327 [59]	Inadequate checks allow users to see the names of tags that were either unpublished or published with restricted view permission.	CWE_200	37 with 22, 23, 24, 25, 50
CVE-2018-17857 [60]	Inadequate checks on the tag search fields can lead to an access level violation.	CWE_863	1 with 22, 23, 24, 25
Password aging problem in Joomla	Joomla does not integrate any mechanism for managing password aging; consequently, users aren't incentivized to update passwords periodically.	CWE_262	2, 3, 5, 6, 7, 8, 11, 12, 15, 17, 20, 22, 23, 24, 25, 26, 27, 28, 29, 30, 66, 110, 144, 146

on each connected component (Section VII), we can exhaustively search the best order of removal steps for each connected component, in order to determine the *optimal solution* for each reduced input set. Thus, we can compare the solutions obtained by different search algorithms for different time budgets to this optimal solution. Note that this approach is tractable for the considered connected components only because they are small, but it is intractable for larger input sets, e.g., the initial or reduced input sets.

2) *Security Vulnerabilities*: The replication package for MST-wi [51] includes 76 metamorphic relations (MRs). These MRs can identify nine vulnerabilities in Jenkins and three vulnerabilities in Joomla using the initial input set, as detailed in Tables II and III, respectively.

For both tables, the first column contains, when available, the CVE identifiers of the considered vulnerabilities. The password aging problem (for both Jenkins and Joomla) and weak password (for Jenkins) are vulnerabilities that were identified during the MST-wi study [6] and therefore do not have CVE identifiers. The second column provides a short description of the vulnerabilities. The third column reports the CWE

ID for each vulnerability. We present two CWE IDs (e.g., CWE 863 and 280) in cases where the CVE report denotes a general vulnerability type (e.g., CWE 863 for incorrect authorization [54]), though a more precise identification (e.g., CWE 280 concerning improper handling of privileges that may result in incorrect authorization) could be applied. Since the 12 considered vulnerabilities are associated to nine different CWE IDs and each vulnerability has a unique CWE ID, we can conclude that the selected subjects cover a diverse set of vulnerability types, thus further improving the generalizability of our results.

The last column in Tables II and III lists identifiers for inputs which were able to trigger the vulnerability using one of the corresponding MRs. For instance, one can detect vulnerability CVE-2018-1999046 in Jenkins by running the MR written for CWE_200 with inputs 2 or 116. For the first two Joomla vulnerabilities, two inputs need to be present at the same time in the input set in order to trigger the vulnerability because, as opposed to most MRs, the corresponding MRs requires two source inputs to generate follow-up inputs. For instance, to detect vulnerability CVE-2018-17857 in Joomla, one needs

input 1 and at least one input amongst inputs 22, 23, 24, or 25.

3) *AIM configurations*: AIM can be configured in different ways to obtain a minimized input set from an initial input set. Such a *configuration* consists in a choice of distance function and algorithm for output clustering (§ VI-B), and a choice of algorithm for action clustering (§ VI-C).

For the sake of conciseness, in Tables V to XI, we denote each configuration by three letters, where L and B respectively denote the Levenshtein and Bag distances, and K, D, and H respectively denote the K-means, DBSCAN, and HDBSCAN clustering algorithms. For instance, BDH denotes that Bag distance and DBSCAN were used for output clustering, and then HDBSCAN for action clustering. These notations are summarized in Table IV.

AIM performs Silhouette analysis (§ VI-A) to determine the hyper-parameters required for these clustering algorithms. We considered the same ranges of values for the hyper-parameters in both output clustering (§ VI-B) and action clustering steps (§ VI-C). For K-means, we select the range $[1, 70]$ for the number of clusters k . In the case of DBSCAN, the range for the distance threshold ϵ is $[2, 10]$ for Jenkins and $[3, 15]$ for Joomla. The range is larger for Joomla because Joomla has a larger number of Web pages than Jenkins. Finally, the range for the minimum number of neighbors n is $[1, 5]$ for both systems. For HDBSCAN, the range for the minimum number n of individuals required to form a cluster is $[2, 8]$ for both systems.

We also determine the hyper-parameters for the genetic search (Section VIII). Related work on whole test suite generation successfully relied on a population of 80 individuals [31]. Since we reduced the problem (Section VII) before applying MOCCO independently to each connected component, which includes fewer inputs than the whole test suite generation [31], we experimented with a lower population size of 20 individuals. Additionally, for RQ1 and RQ2, we set the number of generations for the genetic algorithm to 100, similar to the value considered in previous work [31]. However, for RQ3, to compare the minimized input sets obtained by search algorithms with different time budgets, we use a maximum time budget of 600 seconds as the termination criterion for MOCCO, consistent with MOSA's study [28], while recording the intermediate results over time.

4) *Baselines*: For RQ1 and RQ2, we identify the following baselines against which to compare AIM configurations.

A 2016 survey reported that 57% of metamorphic testing (MT) work used Random Testing (RT) to generate source inputs [5] and, in 2021, 84% of the publications related to MT adopted traditional or improved RT methods to generate source inputs [61]. In the context of test suite minimization, random search is a straightforward baseline against which to compare AIM that is commonly used [21], [22]. This baseline consists in randomly selecting a given number of inputs from the initial input set. This number is determined based on AIM runs. Each AIM run is performed using 18 different configurations (§ X-A3), each leading to a different minimized input set. So, for a fair comparison, we configure random search to select n inputs from the initial input, where n is the size of the largest

TABLE IV
AIM CONFIGURATIONS AND BASELINES FOR RQ1 AND RQ2.

		Output clustering		Action clustering
		Distance	Algorithm	Algorithm
AIM configurations	LKK	Levenshtein	K-means	K-means
	LKD	Levenshtein	K-means	DBSCAN
	LKH	Levenshtein	K-means	HDBSCAN
	LDK	Levenshtein	DBSCAN	K-means
	LDD	Levenshtein	DBSCAN	DBSCAN
	LDH	Levenshtein	DBSCAN	HDBSCAN
	LHK	Levenshtein	HDBSCAN	K-means
	LHD	Levenshtein	HDBSCAN	DBSCAN
	LHH	Levenshtein	HDBSCAN	HDBSCAN
	BKK	Bag	K-means	K-means
	BKD	Bag	K-means	DBSCAN
	BKH	Bag	K-means	HDBSCAN
	BDK	Bag	DBSCAN	K-means
	BDD	Bag	DBSCAN	DBSCAN
	BDH	Bag	DBSCAN	HDBSCAN
	BHK	Bag	HDBSCAN	K-means
	BHD	Bag	HDBSCAN	DBSCAN
	BHH	Bag	HDBSCAN	HDBSCAN
Baselines	R	×	×	Random
	AK	×	×	K-means
	AD	×	×	DBSCAN
	AH	×	×	HDBSCAN

input set produced by the 18 AIM configurations. We repeat this process, for each AIM run, to obtain the same number of input sets for random search as for AIM.

Moreover, Adaptive Random Testing (ART) was proposed to enhance the performance of RT. It is based on the intuition that inputs close to each other are more likely to have similar failure behaviors than inputs further away from each other. Thus, ART generates inputs widely spread across the input domain, in order to find failure with fewer inputs than RT [62]. ART is also commonly used in the context of test suite minimization [16], [38]. It is similar to our action clustering step (§ VI-C), since it is based on partitioning the input space and generating new inputs in blocks that are not already covered [61].

To perform ART, we need to group inputs based on the similarity between their actions. So, we use AIM to perform action clustering directly on the initial input set instead of output classes. Then, for each cluster, we randomly select one input that covers it. Finally, we group the selected inputs to obtain an input set for the ART baseline. This algorithm will stop when we have selected one input from each cluster, and thus, it is not limited by a time budget. Again, we repeat this process for each AIM run. Since we considered the K-means, DBSCAN, and HDBSCAN clustering algorithms, there are three variants of this baseline.

In Tables V to XI, R denotes the random search baseline while AK, AD, and AH denote the ART baselines, using respectively the K-means, DBSCAN, and HDBSCAN clustering algorithms. These notations are summarized in Table IV.

RQ3 aims at determining the contribution of the MOCCO genetic algorithm to input set minimization and comparing it with alternative approaches. We consider random search, a greedy algorithm for the set cover problem (Section VIII), as well as MOSA and NSGA-III (§ II-E2), which are well-known many-objective search algorithms [63]–[66].

We use random search to provide insights on how difficult the problem is and to help assess if search algorithms are

necessary to solve it. Random search starts with an empty input set and, for each iteration, randomly selects one input from the reduced input set (§ X-A1). This process continues until all the objectives are covered.

The greedy algorithm for the set cover problem [46] is one of the best-possible polynomial time approximation algorithm for this problem, with a tight bound on the cost of the solution for the unweighted variant of the problem [45], matching theoretical bounds [67], and that can be efficiently adapted to the weighted variant [46]. Since the weighted set cover problem is close to our problem (Section VIII), we adapt this greedy algorithm to obtain a minimized input set from the reduced input set. The minimized input set is initially empty. For each iteration, the input with the best cost effectiveness is selected, where the cost effectiveness of input in is computed as $\frac{|Cover(in) \cap Uncovered|}{cost(in)}$. Recall that $Cover(in)$ and $cost(in)$ are defined in § III-A, and $Uncovered$ denotes the set of the objectives not yet covered by selected inputs. This process continues until all the objectives are covered.

The initial population for both MOSA and NSGA-III is the reduced input set. For the other parameters, such as population size, mutation rate, and (for NSGA-III) the number of reference points (§ II-E), we use the default values recommended in the original studies [27]–[29]. More precisely, the population size for MOSA is set to 50, while it is automatically computed for NSGA-III based on the number of reference points. We set the crossover rate to 1 for both MOSA and NSGA-III, the same as MOCCO. Finally, we use the same objective functions and mutation operator for both MOSA and NSGA-III, as in MOCCO. Finally, for MOCCO and all baselines, we use the same termination criterion as in MOSA’s study [28], allocating a maximal time budget of 600 seconds for all search algorithms, while recording intermediate results to determine the minimized input sets for smaller time budgets.

5) *Evaluation Metrics*: To reduce the impact of randomness in our experiment, each configuration and baseline was run 50 times on each system, obtaining one minimized input set for each run. Moreover, for the sake of performance analysis, we also recorded the execution time required by AIM to generate minimized input sets. The purpose of the metrics we use for RQ1 and RQ2 is to determine the “best” configuration to run AIM on a given system. But one cannot know, before experimenting with the target system, which configuration would be the “best” for this system. Based on the results from Jenkins and Joomla (see § X-A1), we determine the overall “best” configuration, to be recommended as default for a new system.

For RQ1, we consider the vulnerabilities described in Table II for Jenkins and in Table III for Joomla. We manually investigated the results of the initial input sets to identify the inputs capable of detecting vulnerabilities in the systems under test, so that we can map inputs to vulnerabilities. For each system, we consider that a vulnerability is detected by a minimized input set if it contains at least one input or pair of inputs able to trigger this vulnerability. For the first two Joomla vulnerabilities requiring pairs of inputs, the vulnerability is detected if both inputs are present in the input set. Hence, for each configuration or baseline, our metric is the *vulnerability*

detection rate (VDR), i.e., the total number of vulnerabilities detected by the minimized input sets obtained for the 50 runs, divided by the total number of vulnerabilities detected by the corresponding initial input sets. If VDR is 100%, then we say the configuration or baseline leads to full vulnerability coverage. The overall “best” configuration for Jenkins and Joomla, regarding vulnerability detection, should have a large VDR for both systems, ideally 100%. For each system, we reject configurations and baselines which do not lead to full vulnerability coverage, and then compare the remaining ones to answer RQ2.

RQ2 aims at evaluating the effectiveness of AIM in minimizing the initial input set. Our goal is to identify the AIM configuration generating minimized input sets leading to the minimal execution time for the 76 considered MRs, across the two case studies, and reporting on the execution time saved, compared to executing MST-wi on the full input set. But, to have a fair comparison between MRs execution time obtained respectively with the initial and minimized input sets, we have to take into account the AIM execution time required to minimize the initial input set. Thus, the input set minimization effectiveness is quantified as the sum of AIM execution time to obtain the minimized input set plus MRs execution time with the minimized input set, divided by that of the initial input set. However, since MR execution time is usually large, we cannot collect the time required to execute our 76 MRs on all the input sets generated by all AIM configurations. We estimate it would take thousands of hours for the 1800 runs, resulting from 18 configurations \times 50 repetitions \times 2 case study subjects. For this reason, we rely on three additional metrics, that can be inferred without executing MRs, to identify the “best” configuration. Then, we report on the input set minimization effectiveness obtained by such configuration. Further, to keep the experiment within feasible computation resources, amongst the 50 minimized input sets of the best configuration, we select one with a median cost (§ III-A and V-A), to be representative of the 50 runs.

To determine the “best” AIM configuration, we consider the size of the generated input set (i.e., the number of inputs in it), its cost (§ III-A and V-A), and the time required by the configuration to generate results. Input set size is a direct measure of effectiveness, while cost is an indirect measure, specific to our approach, that is linearly correlated with MR execution time (§ II-B). For these three metrics (size, cost, AIM execution time), we compare, for each system, the AIM configurations and baselines leading to full vulnerability coverage. More precisely, for each metric, we denote by M_i the value of the metric obtained for the i^{th} approach (AIM configurations or baseline); the 50 runs of approach i leading to a sample containing 50 data points.

To compare two samples for M_1 and M_2 , we perform a Mann-Whitney-Wilcoxon test, which is recommended to assess differences in stochastic order for software engineering data analysis [68]. This is a non-parametric test of the null hypothesis that $P(M_1 > M_2) = P(M_1 < M_2)$, i.e., M_1 and M_2 are stochastically equal [69]. Hence, from M_1 and M_2 samples, we obtain the p-value p indicating how likely is the observation of these samples, assuming that M_1 and M_2 are

stochastically equal. If $p \leq 0.05$, we consider it is unlikely that M_1 and M_2 are stochastically equal.

To assess practical significance, we also consider a metric for effect size. An equivalent reformulation of the null hypothesis is $P(M_1 > M_2) + 0.5 \times P(M_1 = M_2) = 0.5$, which can be estimated by counting in the samples the number of times a value for M_1 is larger than a value for M_2 (ties counting for 0.5), then by dividing by the number of comparisons. That way, we obtain the Vargha and Delaney's A_{12} metric [69] which, for the sake of conciseness, we simply denote A in Tables VI to XI. A is considered to be a robust metric for representing effect size in the context of non-parametric methods [70]. A ranges from 0 to 1, where $A = 0$ indicates that $P(M_1 < M_2) = 1$, $A = 0.5$ indicates that $P(M_1 > M_2) = P(M_1 < M_2)$, and $A = 1$ indicates that $P(M_1 > M_2) = 1$.

RQ3 aims at determining the contribution of MOCCO to input set minimization and comparing it with alternative approaches. We use the results of the double-clustering and problem reductions steps to obtain reduced input sets (§ X-A1) from the 50 runs of the “best” configuration. For each reduced input set, each considered algorithm (MOCCO or a baseline) is executed to obtain the corresponding minimized input set. Then, the minimized input sets are checked to determine if the algorithms lead to full vulnerability coverage for this run, and their cost (§ III-A and V-A) is recorded. Again, we denote by M_i the cost for the i^{th} approach. The 50 runs of approach i yield a sample containing 50 data points. As opposed to RQ2, where any run from a configuration/baseline can be compared to any run of another configuration/baseline, we want to compare the cost of the n^{th} minimized input set for an algorithm to the cost of the n^{th} minimized input set of another algorithm, since they are both obtained from the same n^{th} reduced input set. Hence, to compare two samples for M_1 and M_2 , we perform a Wilcoxon signed-rank test, which is a non-parametric paired test [68], with a level of significance of 0.05.

To assess practical significance, we also consider a metric for effect size. Metrics for this test are often defined in terms of the positive-rank sum R^+ and the negative-rank sum R^- [71], [72]. Similarly to the Vargha and Delaney's A_{12} metric [69] we used for the Mann-Whitney-Wilcoxon test to answer RQ2, we use for RQ3 the effect size $E \stackrel{\text{def}}{=} \frac{R^+}{R^+ + R^-}$, so that E ranges from 0 to 1, where $E = 0$ indicates that $M_1 < M_2$ in every case and $E = 1$ indicates that $M_1 > M_2$ in every case.

B. Empirical Results

We first describe the system configurations used to obtain our results (§ X-B1). To answer RQ1, we report the VDR associated with the obtained minimized input sets (§ X-B2). Then, we describe the effectiveness of the input set reduction of the whole AIM approach to answer RQ2 (§ X-B3) and of the MOCCO component to answer RQ3 (§ X-B4).

1) *System Configurations*: We performed all the experiments on a system with the following configurations: a virtual machine installed on professional desktop PCs (Dell G7 7500, RAM 16Gb, Intel(R) Core(TM) i9-10885H CPU @ 2.40GHz)

TABLE V
COVERAGE OF THE JENKINS AND JOOMLA VULNERABILITIES AFTER 50 RUNS OF EACH CONFIGURATION AND BASELINE.

Vulnerability Coverage	System Under Test			
	Jenkins		Joomla	
Configurations or baselines	Nb of detected vulnerabilities	VDR	Nb of detected vulnerabilities	VDR
LKK	450	100.0%	146	97.3%
LKD	371	82.4%	50	33.3%
LKH	379	84.2%	150	100.0%
LDK	450	100.0%	150	100.0%
LDD	400	88.9%	50	33.3%
LDH	400	88.9%	50	33.3%
LHK	450	100.0%	100	66.7%
LHD	403	89.6%	100	66.7%
LHH	447	99.3%	100	66.7%
BKK	450	100.0%	133	88.7%
BKD	403	89.6%	50	33.3%
BKH	410	91.1%	150	100.0%
BDK	450	100.0%	150	100.0%
BDD	338	75.1%	50	33.3%
BDH	450	100.0%	50	33.3%
BHK	450	100.0%	100	66.7%
BHD	404	89.8%	100	66.7%
BHH	428	95.1%	100	66.7%
R	339	75.3%	74	49.3%
AK	447	99.3%	125	83.3%
AD	77	17.1%	22	14.7%
AH	350	77.8%	68	45.3%

and terminal access to a shared remote server with Intel(R) Xeon(R) Gold 6234 CPU (3.30GHz) and 8 CPU cores.

2) *RQ1 - Detected Vulnerabilities*: Results are presented in Table V. Configurations and baselines that lead to full vulnerability coverage for both systems are in green, in yellow if they lead to full vulnerability coverage for one system, and in red if they never lead to full vulnerability coverage. As shown in Table II and Table III, there are 9 vulnerabilities in Jenkins and 3 vulnerabilities in Joomla. Each AIM configuration is executed 50 times to reduce the effect of randomness in our experiments. We consider an AIM configuration to achieve full vulnerability coverage on Jenkins if it achieves $9 \times 50 = 450$ vulnerability detections across all runs. Similarly, full vulnerability coverage on Joomla across all runs is reached if the configuration achieves $3 \times 50 = 150$ vulnerability detections. The execution time of the AIM configurations ranges from 15 to 24 minutes on Jenkins and from 25 to 47 minutes on Joomla. We conclude that such variation across configurations is not significant compared to the time required to execute MRs.

First, note that **the choice of distance function for output clustering does not have a significant impact on vulnerability coverage**. Indeed, apart from LDH and BDH, the results using the Levenshtein or Bag distances are fairly similar (e.g., both LKK and BKK discover 450 vulnerabilities in Jenkins) and seem to only depend on the choice of clustering algorithms. This indicates that the order of words in a Web page is not a relevant distinction when performing clustering for vulnerability coverage. Considering now LDH and BDH, taking into account the order of words can even be detrimental, since they perform equally poorly for Joomla but they differ for Jenkins, where only BDH leads to full vulnerability coverage.

Second, **the choice of clustering algorithm for action clustering seems to be the main factor determining vulner-**

ability coverage. Configurations using DBSCAN as algorithm for the action clustering step never lead to full vulnerability coverage for any system. This indicates that this clustering algorithm poorly fits the data in the input space. This is confirmed by the results obtained for the AD baseline, which only uses DBSCAN on the input space and performs the worst (amongst baselines and AIM configurations) regarding vulnerability coverage. After investigation, the minimized input sets acquired for AD are much smaller compared to those obtained for the other baseline methods. These results cannot be explained by the hyper-parameter as we employed a large range of values (§ X-A3). We conjecture that DBSCAN merges together many action clusters even when the URLs involved in these actions are distinct.

On the other hand, **configurations using K-means for the action clustering step always lead to full vulnerability coverage for Jenkins and lead to the largest vulnerability coverage for Joomla.** This is confirmed by the results obtained for the AK baseline, which only uses K-means on the input space and performs the best (amongst baselines) regarding vulnerability coverage. Indeed, even if this configuration does not lead to full vulnerability coverage, it is very close. In fact, even if it tends to perform worse than AIM configurations that use K-means for action clustering, it tends to perform better than AIM configurations that do not use K-means for action clustering. The success of K-means in achieving better vulnerability coverage on these datasets can be attributed to its ability to handle well-separated clusters. In our case, these clusters are well-separated because of the distinct URLs occurring in the datasets.

Finally, **no baseline reached full vulnerability coverage.** On top of the already mentioned AK and AD baselines, AH performed similarly to random testing (R), indicating that the effect of the HDBSCAN algorithm for action clustering is neutral. The only AIM configuration that performed worse than random testing is BDD, combining DBSCAN (as mentioned before, the worst clustering algorithm regarding vulnerability detection) for both output and action clustering with Bag distance. Only LDK and BDK lead to full vulnerability coverage for both Jenkins and Joomla, and hence are our candidate “best” configurations in terms of VDR. The combination of DBSCAN and K-means was very effective on our dataset since DBSCAN was able to identify dense regions of outputs and K-means allowed for further refinement, forming well-defined action clusters based on URLs.

3) *RQ2 - Input Set Reduction Effectiveness:* To answer RQ2 on the effectiveness of minimization, we compare the input set reduction of baselines and configurations for both Jenkins and Joomla. Amongst them, only the LKK, LDK, LHK, BKK, BDK, BDH, and BHK configurations lead to full vulnerability coverage for Jenkins. Their input set sizes are compared in Table VI, their costs in Table VII, and their AIM execution time in Table VIII. Similarly, only the LKH, LDK, BKH, and BDK configurations lead to full vulnerability coverage for Joomla. Their input set sizes are compared in Table IX, their costs in Table X, and their AIM execution time in Table XI. Configurations with full vulnerability coverage for both Jenkins and Joomla (i.e., LDK and BDK) are in bold.

In these six tables, configurations in each row are compared with configurations in each column. p denotes the statistical significance and A the effect size (§ X-A5). When $p > 0.05$, we consider the metric values obtained from the two configurations not to be significantly different, and hence the cell is left white. Otherwise, the cell is colored, either in green or red. Since we consider input set size and cost and AIM execution time, the smaller the values the better. Thus, green (resp. red) indicates that the configuration in the row is better (resp. worse) than the configuration in the column. The intensity of the color is proportional to the effect size. More precisely, the intensity is $|\delta|$, where $\delta = 2 \times A - 1$ is Cliff’s delta [70]. $|\delta|$ is a number between 0 and 1, where 0 indicates the smallest intensity (the lightest color) and 1 indicates the largest intensity (the darkest color).

For Jenkins, among the candidate best configurations (i.e., LDK and BDK), **BDK performed significantly better than LDK for input set size and cost**, and even if the difference is smaller for AIM execution time, the effect size is also in favor of BDK. As for the other configurations, Table VI on input set sizes and Table VII on input set costs consistently indicate that BDH is the best configuration while LDK is the worst configuration. The other configurations seem equivalent in terms of size. Regarding cost, LKK tends to be the second to last configuration, the other configurations being equivalent. Regarding AIM execution time in Table VIII, the results are more nuanced, BDH is again the best configuration, but this time LKK is the worst configuration instead of LDK. BDH is the only configuration that reached full vulnerability coverage for Jenkins without using the K-means clustering algorithm and it performs significantly better than the other configurations, especially the ones involving K-means for both output and action clustering steps. This indicates, without surprise, that the K-means algorithm takes more resources to be executed. BDH did not lead to full vulnerability coverage for Joomla, so we do not consider it as a candidate for “best” configuration.

For Joomla, **BDK performed significantly better than LDK** for the considered metrics. As for the other configurations, Table IX for input set sizes and Table X for input set costs provide identical results, indicating that LKH and BKH dominate the others while being equivalent. Moreover, BDK dominates LDK, which is the worst configuration. The results are almost identical for AIM execution time in Table XI, with the small difference that LKH performs slightly better than BKH. However, LKH and BKH did not lead to full vulnerability coverage for Jenkins, as opposed to BDK and LDK.

Since we obtained similar results for both Jenkins and Joomla, **we consider BDK to be the “best” AIM configuration.** This is not surprising since Bag distance is less costly to compute than Levenshtein distance (§ VI-B) and we already observed that the order of words in a Web page does not appear to be a relevant distinction for vulnerability coverage (§ X-B2).

As mentioned in § X-B2, no baseline leads to full vulnerability coverage. AD fared poorly and AH performed similarly to random testing R, but AK was much better, with 99.3%

TABLE VI
COMPARISON OF JENKINS INPUT SET SIZES FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

sizes		LKK	LDK	LHK	BKK	BDK	BDH	BHK
LKK	p		$5.4e-15$	$5.1e-1$	$4.2e-1$	$8.8e-1$	$3.1e-20$	$7.2e-1$
	A		0.05	0.54	0.55	0.49	1.0	0.48
LDK	p	$5.4e-15$		$1.8e-17$	$1.4e-15$	$1.8e-14$	$3.2e-20$	$3.8e-14$
	A	0.95		0.99	0.96	0.94	1.0	0.94
LHK	p	$5.1e-1$	$1.8e-17$		$9.8e-1$	$3.4e-1$	$3.0e-20$	$1.6e-1$
	A	0.46	0.01		0.5	0.44	1.0	0.42
BKK	p	$4.2e-1$	$1.4e-15$	$9.8e-1$		$4.1e-1$	$3.2e-20$	$2.1e-1$
	A	0.45	0.04	0.5		0.45	1.0	0.43
BDK	p	$8.8e-1$	$1.8e-14$	$3.4e-1$	$4.1e-1$		$3.2e-20$	$6.9e-1$
	A	0.51	0.06	0.56	0.55		1.0	0.48
BDH	p	$3.1e-20$	$3.2e-20$	$3.0e-20$	$3.2e-20$	$3.2e-20$		$3.2e-20$
	A	0.0	0.0	0.0	0.0	0.0		0.0
BHK	p	$7.2e-1$	$3.8e-14$	$1.6e-1$	$2.1e-1$	$6.9e-1$	$3.2e-20$	
	A	0.52	0.06	0.58	0.57	0.52	1.0	

TABLE VII
COMPARISON OF JENKINS INPUT SET COSTS FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

costs		LKK	LDK	LHK	BKK	BDK	BDH	BHK
LKK	p		$2.5e-16$	$5.9e-5$	$1.5e-2$	$4.4e-3$	$4.1e-18$	$5.4e-3$
	A		0.02	0.73	0.64	0.67	1.0	0.66
LDK	p	$2.5e-16$		$7.0e-18$	$9.5e-18$	$7.0e-18$	$4.1e-18$	$7.0e-18$
	A	0.98		1.0	1.0	1.0	1.0	1.0
LHK	p	$5.9e-5$	$7.0e-18$		$1.0e-1$	$2.4e-1$	$4.1e-18$	$1.4e-1$
	A	0.27	0.0		0.41	0.43	1.0	0.41
BKK	p	$1.5e-2$	$9.5e-18$	$1.0e-1$		$5.8e-1$	$4.1e-18$	$6.1e-1$
	A	0.36	0.0	0.59		0.53	1.0	0.53
BDK	p	$4.4e-3$	$7.0e-18$	$2.4e-1$	$5.8e-1$		$4.1e-18$	$8.2e-1$
	A	0.33	0.0	0.57	0.47		1.0	0.49
BDH	p	$4.1e-18$	$4.1e-18$	$4.1e-18$	$4.1e-18$	$4.1e-18$		$4.1e-18$
	A	0.0	0.0	0.0	0.0	0.0		0.0
BHK	p	$5.4e-3$	$7.0e-18$	$1.4e-1$	$6.1e-1$	$8.2e-1$	$4.1e-18$	
	A	0.34	0.0	0.59	0.47	0.51	1.0	

TABLE VIII
COMPARISON OF JENKINS AIM EXECUTION TIMES FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

times		LKK	LDK	LHK	BKK	BDK	BDH	BHK
LKK	p		$1.5e-6$	$1.0e-11$	$2.0e-2$	$3.8e-5$	$3.1e-18$	$1.3e-9$
	A		0.78	0.89	0.63	0.74	1.0	0.85
LDK	p	$1.5e-6$		$1.2e-5$	$3.9e-3$	$5.4e-1$	$2.6e-18$	$1.1e-3$
	A	0.22		0.75	0.33	0.54	1.0	0.69
LHK	p	$1.0e-11$	$1.2e-5$		$1.7e-9$	$1.4e-3$	$3.2e-17$	$8.7e-1$
	A	0.11	0.25		0.15	0.32	0.98	0.49
BKK	p	$2.0e-2$	$3.9e-3$	$1.7e-9$		$8.0e-3$	$3.0e-18$	$6.7e-7$
	A	0.37	0.67	0.85		0.65	1.0	0.79
BDK	p	$3.8e-5$	$5.4e-1$	$1.4e-3$	$8.0e-3$		$1.5e-17$	$1.1e-2$
	A	0.26	0.46	0.68	0.35		0.99	0.65
BDH	p	$3.1e-18$	$2.6e-18$	$3.2e-17$	$3.0e-18$	$1.5e-17$		$6.3e-16$
	A	0.0	0.0	0.02	0.0	0.01		0.04
BHK	p	$1.3e-9$	$1.1e-3$	$8.7e-1$	$6.7e-7$	$1.1e-2$	$6.3e-16$	
	A	0.15	0.31	0.51	0.21	0.35	0.96	

TABLE IX
COMPARISON OF JOOMLA INPUT SET SIZES FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

sizes		LKH	LDK	BKH	BDK
LKH	p		$4.8e-15$	$1.3e-1$	$2.9e-16$
	A		0.06	0.42	0.06
LDK	p	$4.8e-15$		$1.1e-14$	$4.5e-16$
	A	0.94		0.94	0.94
BKH	p	$1.3e-1$	$1.1e-14$		$7.4e-16$
	A	0.58	0.06		0.06
BDK	p	$2.9e-16$	$4.5e-16$	$7.4e-16$	
	A	0.94	0.06	0.94	

TABLE X
COMPARISON OF JOOMLA INPUT SET COSTS FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

costs		LKH	LDK	BKH	BDK
LKH	p		$1.3e-14$	$6.9e-1$	$3.5e-15$
	A		0.06	0.48	0.05
LDK	p	$1.3e-14$		$1.3e-14$	$7.0e-15$
	A	0.94		0.94	0.94
BKH	p	$6.9e-1$	$1.3e-14$		$3.6e-15$
	A	0.52	0.06		0.05
BDK	p	$3.5e-15$	$7.0e-15$	$3.6e-15$	
	A	0.95	0.06	0.95	

VDR for Jenkins and 83.3% for Joomla. But even if AK had reached full vulnerability coverage for both systems, it would be at a disadvantage compared to AIM configurations. Indeed, over 50 Jenkins runs, the average input set size for

AK was 94.92 inputs, while it ranges from 38 inputs (40% of AK) for BDH to 74.8 inputs (79%) for LDK. The average input set cost for AK was 193,698.94 actions, while it ranges from 70,500.76 actions (36%) for BDH to 152,373.54 actions

TABLE XI

COMPARISON OF JOOMLA AIM EXECUTION TIMES FOR CONFIGURATIONS WITH FULL VULNERABILITY COVERAGE.

times		LKH	LDK	BKH	BDK
LKH	p		$2.5e-18$	$3.6e-3$	$1.2e-18$
	A		0.0	0.33	0.0
LDK	p	$2.5e-18$		$2.8e-18$	$1.5e-18$
	A	1.0		1.0	1.0
BKH	p	$3.6e-3$	$2.8e-18$		$1.3e-18$
	A	0.67	0.0		0.0
BDK	p	$1.2e-18$	$1.5e-18$	$1.3e-18$	
	A	1.0	0.0	1.0	

TABLE XII

COMPARISON OF MRS EXECUTION TIME BEFORE AND AFTER INPUT SET MINIMIZATION. THE PERCENTAGE OF REDUCTION IS ONE MINUS THE RATIO BETWEEN TOTAL EXECUTION TIME AFTER MINIMIZATION AND MRS EXECUTION TIME BEFORE MINIMIZATION.

Execution time (minutes)	Jenkins	Joomla
MRs with initial input set	38,307	20,703
MRs with minimized input set	6119	3675
+ AIM execution time	22	22
= Total execution time	6141	3697
Percentage of Reduction	84%	82%

(79%) for LDK. Over 50 Joomla runs, the average input set size for AK was 70 inputs, while it ranges from 36.02 inputs (51%) for LKH to 41.46 inputs (59%) for LDK. The average input set cost for AK was 2,312,784.58 actions, while it ranges from 580,705.24 actions (25%) for LKH to 872,352.72 actions (38%) for LDK. In short, **all AIM configurations with full vulnerability coverage outperformed the best baseline AK**, which highlights the relevance of our approach in reducing the cost of testing.

Finally, in Table XII, we present the results of executing the MRs using both the initial input set and the minimized input set derived from the best configuration (BDK). In total, by applying AIM, we reduced the execution time of all 76 MRs from 38,307 minutes to 6119 minutes for Jenkins and from 20,703 minutes to 3675 minutes for Joomla, using the minimized input set with median cost. Moreover, executing AIM to obtain this minimized input set required 22 minutes for both systems. Hence, we have a total execution time of 6141 minutes for Jenkins and 3697 minutes for Joomla. As a result, the ratio of the total execution time for the minimized input sets divided by the execution time for the initial input sets is 16.03% for Jenkins and 17.85% for Joomla. In other words, **AIM reduced the execution time by about 84% for Jenkins and more than 82% for Joomla**. This large reduction in execution time demonstrates the effectiveness of our approach in reducing the cost of metamorphic security testing.

4) *RQ3 - Comparison of Search Algorithms*: To answer RQ3, we consider the 50 reduced input sets obtained from the best AIM configuration, namely BDK, and we compare the cost of the minimized input sets obtained by MOCCO and baselines. The cost of the 50 minimized input sets obtained for each genetic algorithm is represented using box plots in Figures 7 and 8 for Jenkins and Joomla, respectively. The results are presented for different time budgets, ranging from 0.2 to 600 seconds, by which time greedy and Many-Objective Coverage and Cost Optimizer have converged. For

TABLE XIII

COMPARISON OF GENETIC ALGORITHMS FOR JENKINS (600 S).

costs		Random	Greedy	MOCCO	MOSA	NSGA-III
Random	p		$1.8e-15$	$1.8e-15$	$1.8e-15$	$1.8e-15$
	E		1.0	1.0	1.0	1.0
Greedy	p	$1.8e-15$		$1.1e-9$	$1.8e-15$	$1.8e-15$
	E	0.0		0.99	0.0	0.0
MOCCO	p	$1.8e-15$	$1.1e-9$		$1.8e-15$	$1.8e-15$
	E	0.0	0.01		0.0	0.0
MOSA	p	$1.8e-15$	$1.8e-15$	$1.8e-15$		$1.8e-15$
	E	0.0	1.0	1.0		1.0
NSGA-III	p	$1.8e-15$	$1.8e-15$	$1.8e-15$	$1.8e-15$	
	E	0.0	1.0	1.0	0.0	

both Jenkins and Joomla, **random search, greedy algorithm, and MOSA quickly converge. Random search and MOSA converge toward sub-optimal solutions**, which is expected since random search is unlikely to determine the best order of removal steps by chance, and MOSA minimizes the cost for each individual objective instead of considering the collective coverage of the selected inputs. **The greedy algorithm finds a good approximation for all runs on both systems**. It finds the optimal solution for most runs (41 out of 50 runs) on Joomla but for only a few runs (6 out of 50 runs) on Jenkins, likely because Jenkins reduced input sets are larger than those of Joomla. **MOCCO finds the optimal solution (§ X-A1) for all 50 runs in 1 second for Jenkins and 0.7 seconds for Joomla**, which is expected since MOCCO is designed to solve this many-objective problem (§ VIII-A). NSGA-III slowly converges for both systems. **NSGA-III does not find the optimal solution within the 600-second time budget for Jenkins, but does so in 600 seconds for Joomla**. This is also expected since NSGA-III explores the entire Pareto front while MOCCO focuses on the region of interest, i.e., around the utopia point of full coverage at no cost. NSGA-III is the only baseline that finds the optimal solution for all runs and within the time budget, but only for one system, **while MOCCO consistently finds the optimal solution in nearly three orders of magnitude faster**.

Furthermore, we conducted statistical tests reported in Table XIII for Jenkins and Tables XIV and XV for Joomla. Since NSGA-III achieves the same result as MOCCO within 600 seconds on Joomla, we also report their results at 400 seconds for a more comprehensive comparison. In these three tables, algorithms in each row are compared with algorithms in each column. p denotes the statistical significance and E the effect size (§ X-A5). As for RQ2, when $p > 0.05$, we consider the costs obtained from the two algorithms not to be significantly different, and hence the cell is left white. Otherwise, the cell is colored, either in green or red. Since we consider input set cost, the smaller the values the better. Thus, green (resp. red) indicates that the algorithm in the row is better (resp. worse) than the algorithm in the column. Finally, as for RQ2, the intensity of the color is computed with $|2 \times E - 1|$.

The results for Jenkins with a time budget of 600 seconds are detailed in Table XIII. The small p -values and effect sizes observed in the MOCCO row indicates that **MOCCO obtained minimized input sets with significantly smaller costs than all the alternative approaches**. Moreover, the

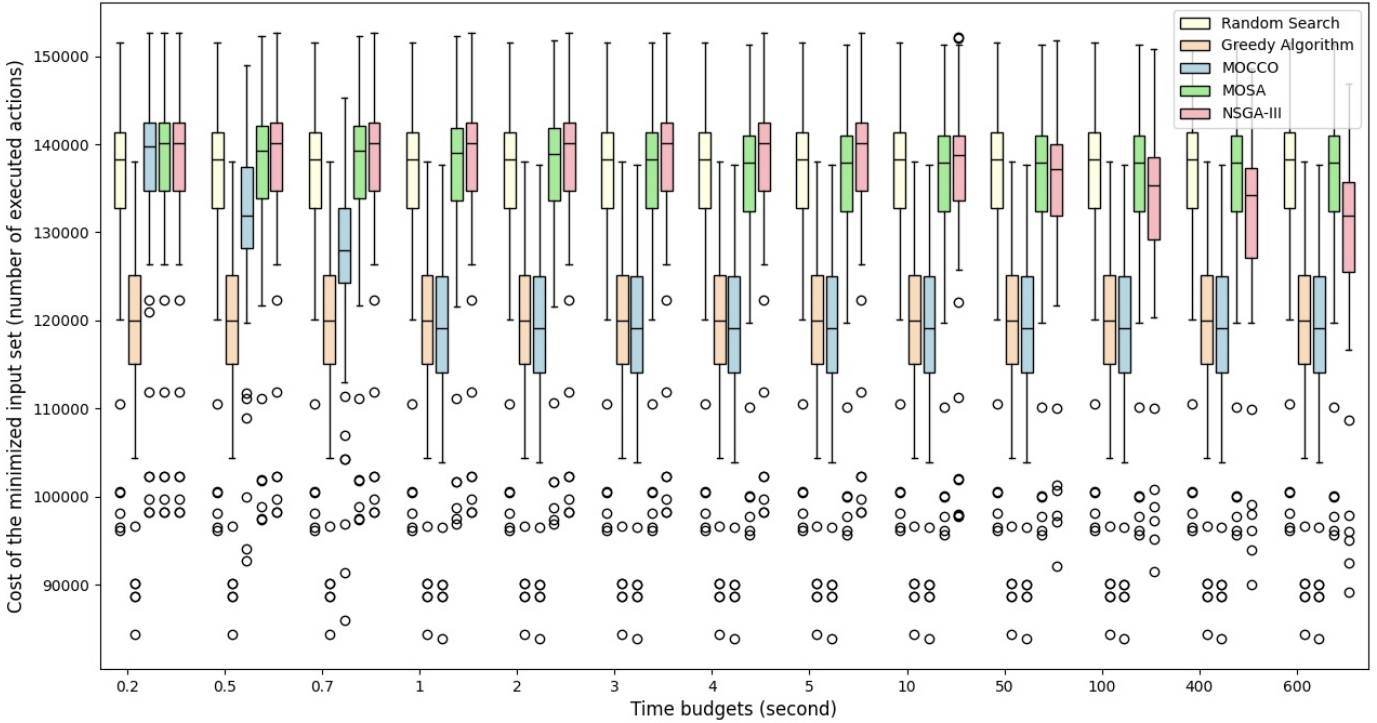


Fig. 7. Jenkins: Cost of the minimized input sets using Random Search, Greedy Algorithm, MOCCO, MOSA, and NSGA-III under different time budgets for 50 runs of BDK.

TABLE XIV
COMPARISON OF GENETIC ALGORITHMS FOR JOOMLA (400 s).

costs		Random	Greedy	MOCCO	MOSA	NSGA-III
Random	p		$1.8e-15$	$1.8e-15$	$1.8e-15$	$4.7e-10$
	E		1.0	1.0	1.0	1.0
Greedy	p	$1.8e-15$		$3.2e-2$	$1.8e-15$	$1.3e-5$
	E	0.0		0.66	0.0	0.15
MOCCO	p	$1.8e-15$	$3.2e-2$		$1.8e-15$	$9.4e-7$
	E	0.0	0.34		0.0	0.11
MOSA	p	$1.8e-15$	$1.8e-15$	$1.8e-15$		$1.8e-5$
	E	0.0	1.0	1.0		0.84
NSGA-III	p	$4.7e-10$	$1.3e-5$	$9.4e-7$	$1.8e-5$	
	E	0.0	0.85	0.89	0.16	

TABLE XV
COMPARISON OF GENETIC ALGORITHMS FOR JOOMLA (600 s).

costs		Random	Greedy	MOCCO	MOSA	NSGA-III
Random	p		$1.8e-15$	$1.8e-15$	$1.8e-15$	$1.8e-15$
	E		1.0	1.0	1.0	1.0
Greedy	p	$1.8e-15$		$3.2e-2$	$1.8e-15$	$3.2e-2$
	E	0.0		0.66	0.0	0.66
MOCCO	p	$1.8e-15$	$3.2e-2$		$1.8e-15$	1.0
	E	0.0	0.34		0.0	0.5
MOSA	p	$1.8e-15$	$1.8e-15$	$1.8e-15$		$1.8e-15$
	E	0.0	1.0	1.0		1.0
NSGA-III	p	$1.8e-15$	$3.2e-2$	1.0	$1.8e-15$	
	E	0.0	0.34	0.5	0.0	

effect size for all baselines but the greedy algorithm is 0, indicating that the minimized input sets obtained by MOCCO consistently have smaller costs across all 50 runs. For the greedy algorithm, the effect size is 0.01 because it performs as well as MOCCO for a few runs (6 out of 50 runs). For

Joomla, the results with a time budget of 400 and 600 seconds are respectively detailed in Tables XIV and XV. For all time budgets, the small p-values and the effect sizes of 0 indicate that **MOCCO performed better than random search and MOSA in all 50 runs**, but the results are more nuanced for the greedy algorithm and NSGA-III. For all time budgets, p-values indicate that **MOCCO performed significantly better than the greedy algorithm**. The latter managed to find the optimal solution for most runs (41 out of 50 runs), yielding an effect size of 0.34, which is still in favor of MOCCO. We conjecture this is because the greedy algorithm does not consider input block as individual objectives and hence, as opposed to MOCCO, does not take into account relevant information when selecting inputs (§ III-B). For a 400-second time budget, the p-value indicates that **MOCCO performed significantly better than NSGA-III**, but NSGA-III managed to find the optimal solution for some runs, hence the effect size of 0.11, which is still in favor of MOCCO. For a 600-second time budget, the p-value is 1.0 and the effect size is 0.5, indicating that both approaches find the same results for all 50 runs, i.e., the optimal solutions. We conjecture that NSGA-III manages to achieve the same results as MOCCO for Joomla but not Jenkins because the former's reduced input set has fewer redundant inputs to be removed compared to Jenkins, making the problem computationally simpler. Since MOCCO was the only search algorithm able to consistently find the optimal solution for Jenkins for every time budget and since, for Joomla, MOCCO finds the optimal solution almost three orders of magnitude faster than the only baseline, i.e., NSGA-III, that manages to obtain the same results, we

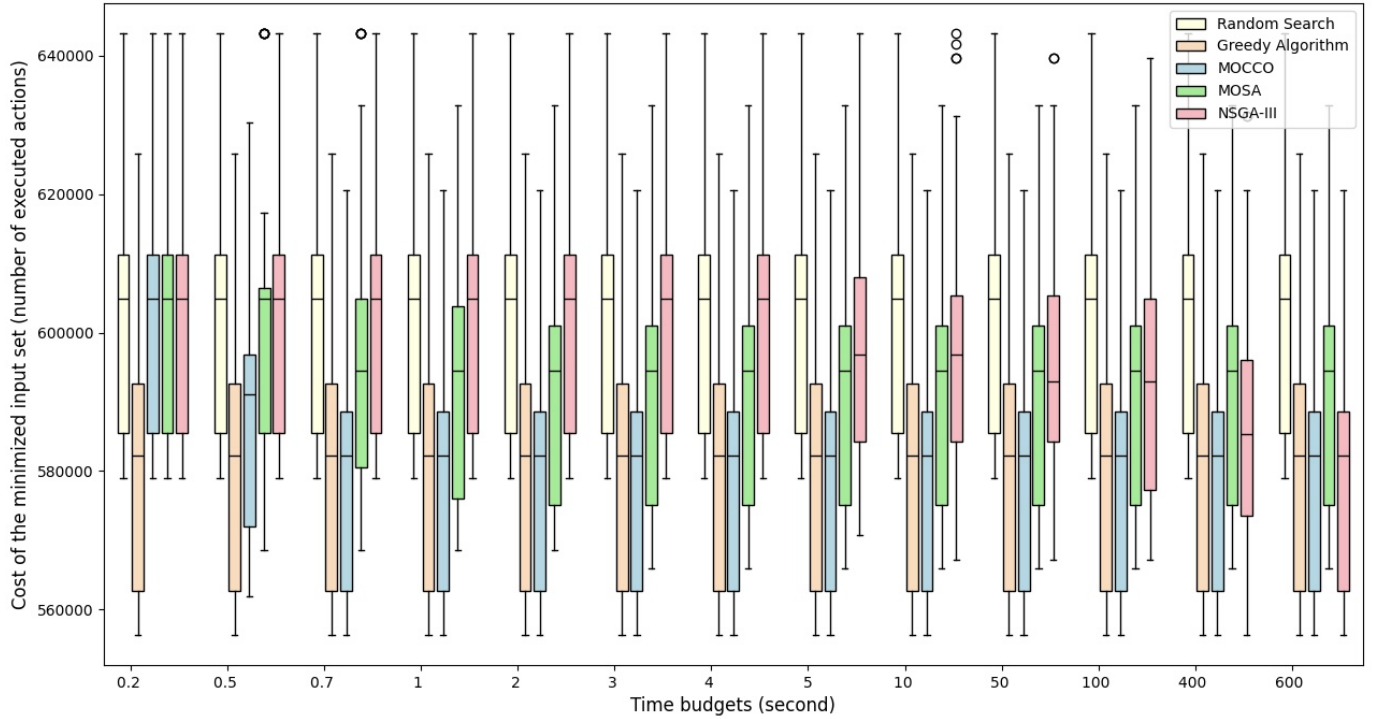


Fig. 8. Joomla: Cost of the minimized input sets using Random Search, Greedy Algorithm, MOCCO, MOSA, and NSGA-III under different time budgets for 50 runs of BDK.

conclude that **MOCCO outperforms all baselines when accounting for both minimized input set cost and execution time.**

Nevertheless, the greedy algorithm finds approximations that are relatively close to those of MOCCO. To determine the significance of the difference in practice, we estimate how this difference in cost translates into a difference in execution time, based on our execution time results (§ X-B3). For Jenkins, it took 6119 minutes to execute a minimized input set of cost 119089 actions. Since the greedy algorithms obtained minimized inputs sets with on average 611.74 more actions and up to 1095 actions, this translates into a difference of 31.43 minutes on average to execute the MRs, up to 56.3 minutes in the worst case. For Joomla, the minimized input set of cost 582181 took 3675 minutes, hence given the difference between the greedy algorithm and MOCCO of 949.32 actions on average with a maximum of 5274 actions, this translates into a difference of 6 minutes on average, up to 33.29 minutes in the worst case. However, relatively to the total execution time, these differences may not have a practical impact. **Therefore, we conclude, from our case studies, that the greedy algorithm, despite its limitations, is a good alternative.**

XI. THREATS TO VALIDITY

In this section, we discuss internal, conclusion, construct, and external validity according to conventional practices [73].

A. Internal Validity

A potential internal threat concerns inadequate data pre-processing, which may adversely impact our results. Indeed,

clustering relies on the computed similarity among the pre-processed outputs and inputs. To address this potential concern, we have conducted a manual investigation of the quality of the clusters obtained without pre-processing. This led us to remove, from the textual content extracted from each Web page, all the content that was shared by many Web pages, like system version, date, or (when present) the menu of the Web page.

For RQ1 on vulnerability detection, one potential threat we face is missing inputs that would be able to exercise a vulnerability or incorrectly considering that an input is able to exercise a vulnerability. To ensure our list of inputs triggering vulnerabilities is complete, one author inspected all the MST-wi execution logs to look for failures.

B. Conclusion Validity

For RQ2, we rely on a non-parametric test (i.e., Mann-Whitney-Wilcoxon test) to evaluate the statistical and practical significance of differences in results, computing p and Vargha and Delaney's A_{12} metric for effect size. Moreover, to deal with the randomness inherent to search algorithms, all the configurations and baselines were executed over 50 runs.

Randomness may also arise from (1) the workload of the machines employed to conduct experiments, potentially slowing down the performance of MST-wi, AIM, and the case study subjects, and (2) the presence of other users interacting with the software under test, which can impact both execution time and system outputs. To address these concerns, we conducted experiments in dedicated environments, ensuring that the study subjects were exclusively utilized by AIM.

C. Construct Validity

The constructs considered in our work are vulnerability detection effectiveness and input set reduction effectiveness. Vulnerability detection effectiveness is measured in terms of vulnerability detection rate. Reduction effectiveness is measured in terms of MR execution time, size and cost of the minimized input set, and AIM execution time for each configuration. As it is expensive to execute all 18 configurations on the MRs, we consider the size of the input set and its cost to select the most efficient configuration. The cost of the input set has been defined in § III-A and shown to be linearly correlated with MR execution time, thus enabling us to evaluate the efficiency of the results.

Finally, we executed the minimized input set obtained from the best configuration on the MRs and compared the obtained execution time, plus the AIM execution time required to minimize the initial input set, with the MRs execution time obtained with the initial input set. Execution time is a direct measure, allowing us to evaluate whether, for systems akin to our case study subjects, AIM should be adopted for making vulnerability testing more efficient and scalable.

D. External Validity

One threat to the generalizability of our results stems from the benchmark that we used. It includes 160 inputs for Jenkins and 148 inputs for Joomla. Furthermore, we considered the list of vulnerabilities in Jenkins and Joomla that were successfully triggered with MST-wi. However, even if in this study we used MST-wi to collect our data, the AIM approach does not depend on a particular data collector, and using or implementing another data collector would enable the use of our approach with other frameworks. Moreover, even if we relied on previously obtained MRs to be sure they detect vulnerabilities in the considered Web systems, AIM is a general approach for metamorphic security testing which does not depend on the considered MRs. Finally, in § X-A1, we highlighted that the different input/output interfaces provided by Jenkins and Joomla, along with the diverse types of vulnerabilities they contain, is in support of the generalizability of our results. Furthermore, the AIM approach can be generalized to other Web systems, if the data collection and pre-processing components are updated accordingly. Nevertheless, further studies involving systems with known vulnerabilities are needed.

XII. RELATED WORK

MT enables the execution of a SUT with a potentially infinite set of inputs thus being more effective than testing techniques requiring the manual specification of either test inputs or oracles. However, in MT, the manual definition of metamorphic relations (MRs) is an expensive activity because it requires that engineers first acquire enough information on the subject under test and then analyze the testing problem to identify MRs. For this reason, in the past, researchers focused on both the definition of methodologies supporting the identification of MRs [74], [75] and the development of techniques for the automated generation of MRs, based on

meta-heuristic search [76], [77] and natural language processing [78], and targeting query-based systems [79] and Cyber-Physical Systems [77].

However, source inputs also impact the effectiveness and performance of MT; indeed, MRs generate follow-up inputs from source inputs and both are executed by the SUT. Consequently, the research community has recently shown increasing interest towards investigating the impact of source inputs on MT. We summarize the most relevant works in the following paragraphs. Note that all these studies focus on general fault detection, while we focus on metamorphic security testing for Web systems. However, our approach could also be applied to fault detection while the approaches below could also be applied to security testing. We therefore compare these approaches without considering their difference in application. However, we excluded from our survey those approaches that study the effect of source and follow-up inputs on the metamorphic testing of systems that largely differ from ours (i.e., sequence alignment programs [80], system validation [81], and deep neural networks [82]). In the following paragraphs, we group the surveyed works into three categories: input generation techniques, input selection techniques, and feedback-directed metamorphic testing.

Input generation techniques for MT use white-box approaches based on knowledge of the source code (mainly, for statement or branch coverage), while we use a black-box approach based on input and output information (Section VI). For instance, a study [83] leveraged the evolutionary search approach EvoSuite [84] to evolve whole input sets in order to obtain inputs that lead to more branch coverage or to different results on the mutated and non-mutated versions of the source code. Another example study [85] leveraged symbolic execution to collect constraints of program branches covered by execution paths, then solved these constraints to generate the corresponding source inputs. Finally, the execution of the generated inputs on the SUT was prioritized based on their contribution regarding uncovered statements. In this case, both generation and prioritization phases were white-box approaches based on branch coverage. Note that, while our approach on input set minimization could be seen as similar to input prioritization, both studies focused on increasing coverage, while we focused on reducing cost while maintaining full coverage.

Input selection techniques share the same objective of our work (i.e., reducing the number of source inputs while maximizing MT effectiveness). Because of its simplicity, random testing (RT) is a common strategy for test suite minimization [21], [22] that has been used in MT [61]. RT was enhanced with Adaptive Random Testing (ART), a technique for obtaining source inputs spread across the input domain with the aim of finding failures with fewer number of inputs than RT. As input selection technique for MT, ART outperforms RT in terms of fault detection [62], which (as in the following studies) was evaluated using the F-measure, i.e., the number of inputs necessary to reveal the first failure. In the AIM approach, our action clustering step (§ VI-C) bears similarities with ART, since we partition inputs based on action parameters which are relevant for our SUT. But, instead of

assuming that close inputs lead to close outputs, we directly used SUT outputs during our output clustering step (§ VI-B), since they are inexpensive to obtain compared to executing MRs. Finally, instead of only counting the number of inputs as in the F-measure (or the size of the input set, in the context of input generation), we considered the cost of each source input as the number of executed actions as surrogate measure, which is tailored to reducing MR execution time in the context of Web systems. Instead of focusing only on distances between source inputs as in ART, another study [61] also investigated distances with follow-up inputs, which is an improvement since usually there are more follow-up inputs than source inputs. This led to the Metamorphic testing-based adaptive random testing (MT-ART) technique, which performed better than other ART algorithms regarding test effectiveness, test efficiency, and test coverage (considering statement, block, and branch coverage). Unfortunately, in the AIM approach, we could not consider follow-up inputs to drive the input selection, since executing MRs to generate these follow-up inputs would defeat our purpose of reducing MR execution time.

Finally, while studies on MT usually focus either on the identification of effective MRs or on input generation/selection, a recent study proposed *feedback-directed metamorphic testing* (FDMT) [86] to determine the next test to perform (both in terms of source input and MR), based on previous test results. They proposed adaptive partition testing (APT) to dynamically select source inputs, based on input categories that lead to fault detection, and a diversity-oriented strategy for MR selection (DOMR) to select an MR generating follow-up inputs that are as different as possible from the already obtained ones. While this approach is promising in general, it is not adapted to our case, where we consider a fixed set of MRs, MR selection being considered outside the scope of this paper. Moreover, since we aim to reduce MR execution time, we cannot execute them and use execution information to guide source input or MR selection during testing. Finally, in our problem definition (Section III), we do not consider source inputs independently from each other, which is why we reduced (Section VII) then minimized (Section VIII) the cost of the input set as a whole.

XIII. CONCLUSION AND FUTURE WORK

As demonstrated in our previous work [6], metamorphic testing alleviates the oracle problem for the automated security testing of Web systems. However, metamorphic testing has shown to be a time-consuming approach. Our approach (AIM) aims to reduce the cost of metamorphic security testing by minimizing the initial input set while preserving its capability at exercising vulnerabilities. Our contributions include 1) a clustering-based black box approach that identifies similar inputs based on their security properties, 2) IMPRO, an approach to reduce the search space as much as possible, then divide it into smaller independent parts, 3) MOCCO, a novel genetic algorithm which is able to efficiently select diverse inputs while minimizing their total cost, and 4) a testing framework automatically performing input set minimization.

We considered 18 different configurations for AIM and we evaluated our approach on two open-source Web systems,

Jenkins and Joomla, in terms of vulnerability detection rate and magnitude of the input set reduction. Our empirical results show that the best configuration for AIM is BDK: Bag distance, DBSCAN to cluster the outputs, and K-means to cluster the inputs. The results show that our approach can automatically reduce MRs execution time by 84% for Jenkins and 82% for Joomla while preserving full vulnerability detection. Across 50 runs, the BDK configuration consistently detected all vulnerabilities in Jenkins and Joomla. We also compared AIM with four baselines common in security testing. Notably, none of the baselines reached full vulnerability coverage. Among them, AK (ART baseline using K-means) emerged as the closest to achieving full vulnerability coverage. All AIM configurations with full vulnerability coverage outperformed this baseline in terms of minimized input set size and cost, demonstrating the effectiveness of our approach in reducing the cost of metamorphic security testing.

Furthermore, we compared the effectiveness of MOCCO, in terms of minimized input set cost and execution time, with four other search algorithms. The results on Jenkins showed that MOCCO obtained minimized input sets with significantly lower costs than all the alternative approaches. The only baseline that could find the optimal solutions on Joomla was NSGA-III, though MOCCO did so almost three orders of magnitude faster. Among the considered alternative search algorithms, greedy was the only algorithm that consistently found results close to those of MOCCO, on both Jenkins and Joomla. Therefore, we conclude that MOCCO outperforms all baselines in terms of minimized input set cost and execution time, while the greedy algorithm, despite its theoretical limitations, remains a viable alternative in practice.

As part of future work, we intend to develop a test case prioritization technique that facilitates earlier vulnerability detection by prioritizing the inputs in the minimized input set that are most likely to detect vulnerabilities.

ACKNOWLEDGMENT

This work is supported by the H2020 COSMOS European project, grant agreement No. 957254, NSERC of Canada under the Discovery and CRC programs, and the Science Foundation Ireland grant 13/RC/2094-2. It is part of a collaborative research program between the University of Ottawa's Nanda laboratory and the SnT centre at the University of Luxembourg.

REFERENCES

- [1] P. X. Mai, F. Pastore, A. Goknil, and L. C. Briand, "A natural language programming approach for requirements-based security testing," *2018 IEEE 29th International Symposium on Software Reliability Engineering (ISSRE)*, pp. 58–69, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:53711718>
- [2] P. X. Mai, A. Goknil, L. K. Shar, F. Pastore, L. C. Briand, and S. Shaame, "Modeling security and privacy requirements: a use case-driven approach," *Information and Software Technology*, vol. 100, pp. 165–182, 2018. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584918300703>
- [3] W. Howden, "Theoretical and empirical studies of program testing," *IEEE Transactions on Software Engineering*, vol. SE-4, no. 4, pp. 293–298, July 1978.

- [4] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The oracle problem in software testing: A survey," *IEEE Transactions on Software Engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [5] S. Segura, G. Fraser, A. B. Sanchez, and A. Ruiz-Cortés, "A survey on metamorphic testing," *IEEE Transactions on Software Engineering*, vol. 42, no. 9, pp. 805–824, Sep. 2016.
- [6] N. C. Bayati, F. Pastore, A. Goknil, and L. C. Briand, "Metamorphic testing for web system security," *IEEE Transactions on Software Engineering*, vol. 49, no. 6, pp. 3430–3471, 2023.
- [7] P. X. Mai, F. Pastore, A. Goknil, and L. C. Briand, "Metamorphic security testing for web systems," *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*, pp. 186–197, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:209202564>
- [8] N. B. Chaleshtari, Y. Marquer, F. Pastore, and L. C. Briand, "Replication package," 2024, <https://zenodo.org/records/13983166>.
- [9] T. Y. Chen, F.-C. Kuo, H. Liu, P.-L. Poon, D. Towey, T. H. Tse, and Z. Q. Zhou, "Metamorphic testing: A review of challenges and opportunities," *ACM Comput. Surv.*, vol. 51, no. 1, jan 2018. [Online]. Available: <https://doi.org/10.1145/3143561>
- [10] OWASP. (2023) Open web application security project. OWASP Foundation. [Online]. Available: <https://www.owasp.org/>
- [11] MITRE. Cwe view: Architectural concepts. MITRE. [Online]. Available: <https://cwe.mitre.org/data/definitions/1008.html>
- [12] —. Cwe-668: Exposure of resource to wrong sphere. MITRE. [Online]. Available: <https://cwe.mitre.org/data/definitions/668.html>
- [13] S. Yoo and M. Harman, "Regression testing minimization, selection and prioritization: a survey," *Softw. Test. Verif. Reliab.*, vol. 22, no. 2, p. 67–120, mar 2012. [Online]. Available: <https://doi.org/10.1002/stv.430>
- [14] B. Miranda and A. Bertolino, "Scope-aided test prioritization, selection and minimization for software reuse," *Journal of Systems and Software*, vol. 131, pp. 528–549, 2017. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121216300875>
- [15] R. Noemmer and R. Haas, "An evaluation of test suite minimization techniques," in *Software Quality: Quality Intelligence in Software and Systems Engineering*, D. Winkler, S. Biffl, D. Mendez, and J. Bergs-mann, Eds. Cham: Springer International Publishing, 2020, pp. 51–66.
- [16] E. Cruciani, B. Miranda, R. Verdecchia, and A. Bertolino, "Scalable approaches for test suite reduction," in *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, May 2019, pp. 419–429.
- [17] R. Pan, T. A. Ghaleb, and L. Briand, "Atm: Black-box test case minimization based on test code similarity and evolutionary search," in *Proceedings of the 45th International Conference on Software Engineering*, ser. ICSE '23. IEEE Press, 2023, p. 1700–1711. [Online]. Available: <https://doi.org/10.1109/ICSE48619.2023.00146>
- [18] C. Arora, M. Sabetzadeh, L. Briand, and F. Zimmer, "Automated extraction and clustering of requirements glossary terms," *IEEE Transactions on Software Engineering*, vol. 43, no. 10, pp. 918–945, 2017.
- [19] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, ser. KDD'96. Portland, Oregon: AAAI Press, 1996, p. 226–231.
- [20] L. McInnes, J. Healy, and S. Astels, "hdbscan: Hierarchical density based clustering," *Journal of Open Source Software*, vol. 2, no. 11, p. 205, 2017. [Online]. Available: <https://doi.org/10.21105/joss.00205>
- [21] S. Wang, S. Ali, and A. Gotlieb, "Cost-effective test suite minimization in product lines using search techniques," *Journal of Systems and Software*, vol. 103, pp. 370–391, 2015. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121214001757>
- [22] M. Zhang, S. Ali, and T. Yue, "Uncertainty-wise test case generation and minimization for cyber-physical systems," *Journal of Systems and Software*, vol. 153, pp. 1–21, 2019. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121219300561>
- [23] B. Li, J. Li, K. Tang, and X. Yao, "Many-objective evolutionary algorithms: A survey," *ACM Comput. Surv.*, vol. 48, no. 1, sep 2015. [Online]. Available: <https://doi.org/10.1145/2792984>
- [24] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: Nsga-ii," *IEEE Transactions on Evolutionary Computation*, vol. 6, no. 2, pp. 182–197, April 2002.
- [25] E. Zitzler, M. Laumanns, and L. Thiele, "Spea2: Improving the strength pareto evolutionary algorithm," Computer Engineering and Communication Networks Lab (TIK), Swiss Federal Institute of Technology (ETH), Zurich, Tech. Rep. 103, 2001.
- [26] M. Kim, T. Hiroyasu, M. Miki, and S. Watanabe, "Spea2+: Improving the performance of the strength pareto evolutionary algorithm 2," in *Parallel Problem Solving from Nature - PPSN VIII*, X. Yao, E. K. Burke, J. A. Lozano, J. Smith, J. J. Merelo-Guervós, J. A. Bullinaria, J. E. Rowe, P. Tiño, A. Kabán, and H.-P. Schwefel, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 742–751.
- [27] K. Deb and H. Jain, "An evolutionary many-objective optimization algorithm using reference-point-based nondominated sorting approach, part i: Solving problems with box constraints," *IEEE Transactions on Evolutionary Computation*, vol. 18, no. 4, pp. 577–601, Aug 2014.
- [28] A. Panichella, F. M. Kifetew, and P. Tonella, "Reformulating branch coverage as a many-objective optimization problem," in *2015 IEEE 8th International Conference on Software Testing, Verification and Validation (ICST)*. Graz, Austria: IEEE, April 2015, pp. 1–10.
- [29] H. Jain and K. Deb, "An evolutionary many-objective optimization algorithm using reference-point based nondominated sorting approach, part ii: Handling constraints and extending to an adaptive approach," *IEEE Transactions on Evolutionary Computation*, vol. 18, no. 4, pp. 602–622, Aug 2014.
- [30] P. Ammann and J. Offutt, *Introduction to Software Testing*. Cambridge: Cambridge University Press, 2016.
- [31] G. Fraser and A. Arcuri, "Whole test suite generation," *IEEE Transactions on Software Engineering*, vol. 39, no. 2, pp. 276–291, Feb 2013.
- [32] A. Arcuri, "It does matter how you normalise the branch distance in search based software testing," in *Third International Conference on Software Testing, Verification and Validation*. Paris, France: IEEE, April 2010, pp. 205–214.
- [33] Z. Li, M. Harman, and R. M. Hierons, "Search algorithms for regression test case prioritization," *IEEE Transactions on Software Engineering*, vol. 33, no. 4, pp. 225–237, April 2007.
- [34] S. Busygin, O. Prokopyev, and P. M. Pardalos, "Biclustering in data mining," *Computers & Operations Research*, vol. 35, no. 9, pp. 2964–2987, 2008, part Special Issue: Bio-inspired Methods in Combinatorial Optimization. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0305054807000159>
- [35] B. Pontes, R. Giráldez, and J. S. Aguilar-Ruiz, "Biclustering on expression data: A review," *Journal of Biomedical Informatics*, vol. 57, pp. 163–180, 2015. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1532046415001380>
- [36] M. Attaoui, H. Fahmy, F. Pastore, and L. Briand, "Black-box safety analysis and retraining of dnns based on feature extraction and clustering," *ACM Trans. Softw. Eng. Methodol.*, vol. 32, no. 3, apr 2023. [Online]. Available: <https://doi.org/10.1145/3550271>
- [37] M. O. Attaoui, H. Fahmy, F. Pastore, and L. Briand, "Dnn explanation for safety analysis: an empirical evaluation of clustering-based approaches," *arXiv*, 2023.
- [38] H. Hemmati, A. Arcuri, and L. Briand, "Achieving scalable model-based testing through test case diversity," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 22, no. 1, pp. 1–42, 2013.
- [39] J. Thomé, L. K. Shar, D. Bianculli, and L. Briand, "Search-driven string constraint solving for vulnerability detection," in *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. Buenos Aires, Argentina: IEEE, 2017, pp. 198–208.
- [40] S. Zhang, Y. Hu, and G. Bian, "Research on string similarity algorithm based on levenshtein distance," in *2017 IEEE 2nd Advanced Information Technology, Electronic and Automation Control Conference (IAEAC)*. Chongqing, China: IEEE, 2017, pp. 2247–2251.
- [41] S. Mergen, "Extending the bag distance for string similarity search," *SN Comput. Sci.*, vol. 4, no. 2, dec 2022. [Online]. Available: <https://doi.org/10.1007/s42979-022-01502-5>
- [42] I. Bartolini, P. Ciaccia, and M. Patella, "String matching with metric trees using an approximate distance," in *String Processing and Information Retrieval: 9th International Symposium, SPIRE 2002 Lisbon, Portugal, September 11–13, 2002 Proceedings 9*. Berlin, Heidelberg: Springer, 2002, pp. 271–283.
- [43] M. Biagiola, A. Stocco, F. Ricca, and P. Tonella, "Diversity-based web test generation," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2019. New York, NY, USA: Association for Computing Machinery, 2019, p. 142–153. [Online]. Available: <https://doi.org/10.1145/3338906.3338970>
- [44] B. Korte and J. Vygen, *Combinatorial Optimization: Theory and Algorithms*, 6th ed. Springer Publishing Company, Incorporated, 2018.
- [45] P. Slavík, "A tight analysis of the greedy algorithm for set cover," in *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, ser. STOC '96. New York, NY, USA: Association for Computing Machinery, 1996, p. 435–441. [Online]. Available: <https://doi.org/10.1145/237814.237991>
- [46] V. V. Vazirani, *Approximation algorithms*. Berlin, Heidelberg: Springer-Verlag, 2001.

- [47] B. Korte and R. Schrader, "On the existence of fast approximation schemes," in *Nonlinear Programming 4*, O. L. Mangasarian, R. R. Meyer, and S. M. Robinson, Eds. Madison, Wisconsin: Academic Press, 1981, pp. 415–437. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780124686625500203>
- [48] Eclipse Foundation, "Jenkins ci/cd server," <https://jenkins.io/>, 2018.
- [49] Joomla!, "Joomla," <https://www.joomla.org/>, 2018.
- [50] (2018) Selenium Web Testing Framework, <https://www.seleniumhq.org/>. Selenium.
- [51] N. B. Chaleshtari, F. Pastore, A. Goknil, and L. C. Briand, "Replication package," 2023, <https://doi.org/10.5281/zenodo.7702754>.
- [52] MITRE. (2018, Nov.) Cve-2018-1000406, concerns cwe-276. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1000406>
- [53] —. (2018) Cve-2018-1000409, concerns otc-sess-003. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1999003>
- [54] —. (2018, Nov.) Cve-2018-1999003, concerns otc-authz-002. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1999003>
- [55] —. (2018, Nov.) Cve-2018-1999004, concerns otc-authz-002. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1999004>
- [56] —. (2018) Cve-2018-1999006, concerns cwe-138. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1999006>
- [57] —. (2018, Nov.) Cve-2018-1999046, concerns otc-authz-002. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-1999046>
- [58] —. (2020) Cve-2020-2162, concerns otc-inpval-003. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2020-2162>
- [59] —. (2018) Cve-2018-11327, concerns cwe-200. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-11327>
- [60] —. (2018) Cve-2018-17857, concerns cwe-200. MITRE. [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2018-17857>
- [61] Z.-w. Hui, X. Wang, S. Huang, and S. Yang, "MT-ART: A Test Case Generation Method Based on Adaptive Random Testing and Metamorphic Relation," *IEEE Transactions on Reliability*, vol. 70, no. 4, pp. 1397–1421, Dec 2021.
- [62] A. C. Barus, T. Y. Chen, F.-C. Kuo, H. Liu, and H. W. Schmidt, "The impact of source test case selection on the effectiveness of metamorphic testing," in *Proceedings of the 1st International Workshop on Metamorphic Testing*, ser. MET '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 5–11. [Online]. Available: <https://doi.org/10.1145/2896971.2896977>
- [63] N. Gupta, A. Sharma, and M. K. Pachariya, "Multi-objective test suite optimization for detection and localization of software faults," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 6, Part A, pp. 2897–2909, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1319157819313850>
- [64] A. Kiran, W. H. Butt, M. W. Anwar, F. Azam, and B. Maqbool, "A comprehensive investigation of modern test suite optimization trends, tools and techniques," *IEEE Access*, vol. 7, pp. 89 093–89 117, 2019.
- [65] M. Naz, Z. Anwaar, and W. H. Butt, "Automated white box test case generation for statement coverage using u-nsga-iii," in *2023 17th International Conference on Open Source Systems and Technologies (ICOSST)*, 2023, pp. 1–6.
- [66] Z. Aghababaeian, M. Abdellatif, M. Dadkhah, and L. Briand, "Deepgd: A multi-objective black-box test selection approach for deep neural networks," *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 6, jun 2024. [Online]. Available: <https://doi.org/10.1145/3644388>
- [67] I. Dinur and D. Steurer, "Analytical approach to parallel repetition," in *Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing*, ser. STOC '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 624–633. [Online]. Available: <https://doi-org.proxy.bnl.lu/10.1145/2591796.2591884>
- [68] A. Arcuri and L. Briand, "A hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering," *Softw. Test. Verif. Reliab.*, vol. 24, no. 3, p. 219–250, may 2014. [Online]. Available: <https://doi-org.proxy.bnl.lu/10.1002/stvr.1486>
- [69] A. Vargha and H. D. Delaney, "A critique and improvement of the cl common language effect size statistics of mcgraw and wong," *Journal of Educational and Behavioral Statistics*, vol. 25, no. 2, pp. 101–132, 2000. [Online]. Available: <https://doi.org/10.3102/10769986025002101>
- [70] B. Kitchenham, L. Madeyski, D. Budgen, J. Keung, P. Brereton, S. Charters, S. Gibbs, and A. Pohthong, "Robust statistical methods for empirical software engineering," *Empirical Softw. Engg.*, vol. 22, no. 2, p. 579–630, apr 2017. [Online]. Available: <https://doi.org/10.1007/s10664-016-9437-5>
- [71] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *J. Mach. Learn. Res.*, vol. 7, p. 1–30, dec 2006.
- [72] D. S. Kerby, "The simple difference formula: An approach to teaching nonparametric correlation," *Comprehensive Psychology*, vol. 3, p. 11.IT.3.1, 2014. [Online]. Available: <https://doi.org/10.2466/11.IT.3.1>
- [73] C. Wohlin, P. Runeson, M. Hst, M. C. Ohlsson, B. Regnell, and A. Wessln, *Experimentation in Software Engineering*. Heidelberg, Germany: Springer Publishing Company, Incorporated, 2012.
- [74] T. Y. Chen, P.-L. Poon, and X. Xie, "Metric: Metamorphic relation identification based on the category-choice framework," *Journal of Systems and Software*, vol. 116, pp. 177–190, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121215001624>
- [75] C.-A. Sun, A. Fu, P.-L. Poon, X. Xie, H. Liu, and T. Y. Chen, "Metric⁺⁺: A metamorphic relation identification technique based on input plus output domains," *IEEE Transactions on Software Engineering*, vol. 47, no. 9, pp. 1764–1785, 2021.
- [76] B. Zhang, H. Zhang, J. Chen, D. Hao, and P. Moscato, "Automatic discovery and cleansing of numerical metamorphic relations," in *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. Cleveland, USA: IEEE, 2019, pp. 235–245.
- [77] J. Ayerdi, V. Terragni, A. Arrieta, P. Tonella, G. Sagardui, and M. Arratibel, "Generating metamorphic relations for cyber-physical systems with genetic programming: An industrial case study," in *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2021. New York, NY, USA: Association for Computing Machinery, 2021, p. 1264–1274. [Online]. Available: <https://doi.org/10.1145/3468264.3473920>
- [78] A. Blasi, A. Gorla, M. D. Ernst, M. Pezzè, and A. Carzaniga, "Memo: Automatically identifying metamorphic relations in javadoc comments for test automation," *Journal of Systems and Software*, vol. 181, p. 111041, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121221001382>
- [79] S. Segura, J. C. Alonso, A. Martín-López, A. Durán, J. Troya, and A. Ruiz-Cortés, "Automated generation of metamorphic relations for query-based systems," in *Proceedings of the 7th International Workshop on Metamorphic Testing*, ser. MET '22. New York, NY, USA: Association for Computing Machinery, 2023, p. 48–55. [Online]. Available: <https://doi.org/10.1145/3524846.3527338>
- [80] J. Y. Tang, A. Yang, T. Y. Chen, and J. W. Ho, "Harnessing multiple source test cases in metamorphic testing: A case study in bioinformatics," in *2017 IEEE/ACM 2nd International Workshop on Metamorphic Testing (MET)*. Buenos Aires, Argentina: IEEE, 2017, pp. 10–13.
- [81] M. Zhang, J. W. Keung, T. Y. Chen, and Y. Xiao, "Validating class integration test order generation systems with metamorphic testing," *Information and Software Technology*, vol. 132, p. 106507, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584920302470>
- [82] J. Zhou, K. Qiu, Z. Zheng, T. Y. Chen, and P.-L. Poon, "Using metamorphic testing to evaluate dnn coverage criteria," in *2020 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)*. Coimbra, Portugal: IEEE, 2020, pp. 147–148.
- [83] P. Saha and U. Kanewala, "Fault detection effectiveness of source test case generation strategies for metamorphic testing," in *Proceedings of the 3rd International Workshop on Metamorphic Testing*, ser. MET '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 2–9. [Online]. Available: <https://doi.org/10.1145/3193977.3193982>
- [84] G. Fraser and A. Arcuri, "Evosuite: Automatic test suite generation for object-oriented software," in *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering*, ser. ESEC/FSE '11. New York, NY, USA: Association for Computing Machinery, 2011, p. 416–419. [Online]. Available: <https://doi-org.proxy.bnl.lu/10.1145/2025113.2025179>
- [85] C.-a. Sun, B. Liu, A. Fu, Y. Liu, and H. Liu, "Path-directed source test case generation and prioritization in metamorphic testing," *Journal of Systems and Software*, vol. 183, p. 111091, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121221001886>
- [86] C.-A. Sun, H. Dai, H. Liu, and T. Y. Chen, "Feedback-directed metamorphic testing," *ACM Trans. Softw. Eng. Methodol.*, vol. 32, no. 1, feb 2023. [Online]. Available: <https://doi.org/10.1145/3533314>



Nazanin Bayati Chaleshtari is a member of Nanda Lab and is currently working toward the PhD degree in the School of EECS, University of Ottawa. She gained valuable practical experience from her work with BlackBerry's security research and development lab in Canada. Throughout her academic career, she has been the recipient of several academic awards, including a PhD admission scholarship, an international doctoral scholarship from the University of Ottawa, and an honourable award for being an outstanding student during her master's degree

at the Iran University of Science and Technology. She was also ranked the best student among all computer engineering students at the Iran University of Science and Technology in 2019. Her research interests include automated software testing concerning security testing, applied data science and empirical software engineering.



Lionel C. Briand is professor of software engineering and has shared appointments between (1) The University of Ottawa, Canada, and (2) The Lero SFI Centre—the national Irish centre for software research—hosted by the University of Limerick, Ireland. In collaboration with colleagues, for over 30 years, he has run many collaborative research projects with companies in the automotive, satellite, aerospace, energy, financial, and legal domains. Lionel has held various engineering, academic, and leading positions in seven countries. He currently

holds a Canada Research Chair (Tier 1) on "Intelligent Software Dependability and Compliance" and is the director of Lero, the national Irish centre for software research. Lionel was elevated to the grades of IEEE Fellow and ACM Fellow for his work on software testing and verification. Further, he was granted the IEEE Computer Society Harlan Mills award, the ACM SIGSOFT outstanding research award, and the IEEE Reliability Society engineer-of-the-year award. He also received an ERC Advanced grant in 2016 on modelling and testing cyber-physical systems, the most prestigious individual research award in the European Union and was elected a fellow of the Academy of Science, Royal Society of Canada in 2023. His research interests include: Testing and verification, trustworthy AI, search-based software engineering, model-driven development, requirements engineering, and empirical software engineering. More details can be found at: <http://www.lbriand.info>.



Yoann Marquer is a post-doctoral researcher at the Interdisciplinary Centre for Security, Reliability and Trust (SnT). He obtained a ministerial scholarship for his PhD in Computer Science, obtained in 2015 from the University of Paris-Est Créteil, after working on the algorithmic completeness and implicit complexity of imperative programming languages. He then worked with academic (universities and the Inria national institute) and industrial partners in several security-related, EU-founded projects. His research interests concern non-functional properties

of computation, especially security, including novel security metrics and countermeasures as well as source code testing and analysis to detect security vulnerabilities and refactoring to make it more secure.



Fabrizio Pastore is Chief Scientist II at the Interdisciplinary Centre for Security, Reliability and Trust (SnT), University of Luxembourg. He obtained his PhD in Computer Science in 2010 from the University of Milano - Bicocca. His research interests concern automated software testing, including security testing and testing of AI-based systems; his work relies on the integrated analysis of different types of artefacts (e.g., requirements, models, source code, and execution traces). He is active in several industry partnerships and national, ESA, and EU-

funded research projects.

AIM: Automated Input Set Minimization for Metamorphic Security Testing

Nazanin Bayati Chaleshtari, Yoann Marquer, Fabrizio Pastore, *Member, IEEE*, and Lionel C. Briand, *Fellow, IEEE*

APPENDIX

We prove in this appendix theoretical results mentioned in the main paper and which demonstrates the correctness of the AIM approach. This includes Theorem 1 (Appendix C), presented in Section III-C on the objective functions, Proposition 3 and Theorem 2 (Appendix D) presented in Section VII-E on local dominance, Proposition 8 and Theorem 3 (Appendix E) presented in Section VII-F on dividing the problem, and desirable properties of roofer (Theorem 4) and miser (Theorem 5) populations (Appendix F) presented in Section VIII-F on the MOCCO population update. Appendices A and B present intermediate results.

A. Input Coverage and Cost

Lemma 1 (Non-Empty Coverage). *For every input in , we have $Cover(in) \neq \emptyset$.*

Proof. We consider only inputs containing at least one action (§VI-B). Let $act_1 = action(in, 1)$ be the first action of in . act_1 has an output in $outCl_1 = OutputClass(in, 1)$ (§VI-B). act_1 is in the action set $actSet_1 = ActionSet(outCl_1)$ (§VI-C). After clustering of this action set, let $bl_1 = Subclass(in, 1) = ActionSubclass(act_1, actSet_1)$ be the input block of the action act_1 executed in in (§VI-C). Therefore, we have $bl_1 \in Cover(in)$. \square

Lemma 2 (Cost is Positive). *For every input in , we have $cost(in) > 0$. For every input set I , we have $cost(I) \geq 0$, and $cost(I) = 0$ if and only if $I = \emptyset$.*

Proof. Since we removed inputs with no cost (§III-A), for each input in , $cost(in) > 0$. The cost of an input set is the sum of the cost of its inputs, hence $cost(I) \geq 0$. In particular, if $I = \emptyset$, then $cost(I) = 0$. Finally, because the cost is positive, the only case where $cost(I) = 0$ is when $I = \emptyset$. \square

B. Redundancy

In this section, we prove that our characterization of redundancy is sound regarding the coverage of an input set

N. Bayati Chaleshtari and L. Briand are with the School of Electrical and Computer Engineering of University of Ottawa, Canada, Y. Marquer and F. Pastore are with the Interdisciplinary Centre for Security, Reliability, and Trust (SnT) of the University of Luxembourg, Luxembourg, and L. Briand is also with Lero SFI Centre for Software Research and University of Limerick, Ireland. Part of this work was done when L. Briand was affiliated with the Interdisciplinary Centre for Security, Reliability, and Trust (SnT) of the University of Luxembourg.

E-mail: n.bayati@uottawa.ca, yoann.marquer@uni.lu, fabrizio.pastore@uni.lu, lbriand@uottawa.ca

Manuscript received Month DD, 2024; revised Month DD, 2024.

(Proposition 1), then we introduce Lemma 4 which is useful to prove Lemma 6 in Appendix C and Proposition 8 in Appendix E.

First, if an input in is in the considered input set I , then its redundancy is not negative. Indeed, there is always at least one input in I that covers the input blocks covered by in , which is in itself.

Lemma 3 (Redundancy is Non-negative). *Let I be an input set and in be an input. If $in \in I$, then $redundancy(in, I) \geq 0$.*

Proof. Let in be an input in I . Let $bl \in Cover(in)$ be any block covered by in . Because $bl \in Cover(in)$, we have that $in \in Inputs(bl)$ (§III-A). Moreover, $in \in I$, so $in \in Inputs(bl) \cap I$, thus we have $superpos(bl, I) \geq 1$ (§III-C). Therefore, for any $bl \in Cover(in)$ we have $superpos(bl, I) \geq 1$. So, $\min\{superpos(bl, I) \mid bl \in Cover(in)\} \geq 1$. By subtracting 1, we have $redundancy(in, I) \geq 0$ (§III-C). \square

We prove that our characterization of redundancy is sound regarding the coverage of an input set. In other words, if an input is redundant in an input set, then it can be removed without reducing the coverage of the input set:

Proposition 1 (Redundancy Soundness). *Let $I \subseteq I_{init}$ be an input set and $in \in I$. If $in \in Redundant(I)$, then $Cover(I \setminus \{in\}) = Cover(I)$.*

Proof. $Cover(I) = Cover(I \setminus \{in\}) \cup Cover(in)$ (§VI-C). We assume that in is redundant in I . So, according to Lemma 3, we have $redundancy(in, I) > 0$ (§III-C). Thus, for each $bl \in Cover(in)$, we have that $superpos(bl, I) \geq 2$. Hence, according to our definition of superposition (§III-C), there are at least two inputs in_1 and in_2 in I which also belong to $Inputs(bl)$. Because we have $bl \in Cover(in)$, we have that $in \in Inputs(bl)$ (§III-A). So, in is one of these inputs. We assume this is the first one i.e., $in_1 = in$. Therefore, for each $bl \in Cover(in)$, there exists an input $in_2 \in Inputs(bl) \cap I$ which is not in . We denote $in_2 = in_{bl}$ this input. For each $bl \in Cover(in)$, we have $in_{bl} \in Inputs(bl)$. So, we have $bl \in Cover(in_{bl})$ (§III-A). Moreover, $in_{bl} \in I$ and is not in , so is in $I \setminus \{in\}$. Thus (§VI-C) we have $bl \in Cover(I \setminus \{in\})$. Therefore, $Cover(in) \subseteq Cover(I \setminus \{in\})$. Finally, $Cover(I) = Cover(I \setminus \{in\}) \cup Cover(in) = Cover(I \setminus \{in\})$. \square

Finally, removing an input in an input set may update the redundancy of other inputs:

Lemma 4 (Removing an Input). *Let I be an input set and $in_2, in_1 \in I$ be two inputs in this input set.*

$$\begin{aligned} \text{redundancy}(in_1, I) - 1 &\leq \text{redundancy}(in_1, I \setminus \{in_2\}) \\ &\leq \text{redundancy}(in_1, I) \end{aligned}$$

Proof. Let $bl \in \text{Cover}(in_1)$. We denote $c_1 = |\text{Inputs}(bl) \cap I|$ and $c_2 = |\text{Inputs}(bl) \cap (I \setminus \{in_2\})|$. If $in_2 \in \text{Inputs}(bl)$ then $c_2 = c_1 - 1$. Otherwise, $c_2 = c_1$. Therefore, for every $bl \in \text{Cover}(in_1)$ we have $c_1 - 1 \leq c_2 \leq c_1$. This includes the input blocks minimizing the cardinality in the definition of redundancy (§III-C). Hence the result. \square

C. Valid Orders of Removal Steps

In this section, we first prove Lemma 5 that establishes that orders of removal steps contain inputs without repetition. Then, we introduce subsequences and prove Lemma 7, used in the rest of the section and in proofs of Appendix E. Finally, we prove Theorem 1 presented in §III-C.

Lemma 5 (Order without Repetition). *Let I be an input set. If $[in_1, \dots, in_n] \in \text{ValidOrders}(I)$, then in_1, \dots, in_n are distinct.*

Proof. The proof is done by induction on n .

If $n = 0$, then $[in_1, \dots, in_n] = []$ is empty, hence there are no two identical inputs.

We now consider $[in_1, \dots, in_n, in_{n+1}] \in \text{ValidOrders}(I)$. By induction hypothesis, in_1, \dots, in_n are distinct. According to the definition of valid order of removal steps (§III-C), we have $[in_1, \dots, in_n, in_{n+1}]$ only if $in_{n+1} \in \text{Redundant}(I \setminus \{in_1, \dots, in_n\})$.

Since $\text{Redundant}(I) = \{in \in I \mid \text{redundancy}(in, I) > 0\}$, we have $in_{n+1} \in I \setminus \{in_1, \dots, in_n\}$.

Therefore, in_{n+1} is distinct from the previous inputs in_1, \dots, in_n . Hence the result, which concludes the induction step. \square

Lemma 6 (Redundant Inputs After Reduction). *Let I be an input set. For each subset of inputs $\{in_1, \dots, in_n\} \subseteq I$, we have $\text{Redundant}(I \setminus \{in_1, \dots, in_n\}) \subseteq \text{Redundant}(I)$.*

Proof. The proof is done by induction on n .

If $n = 0$ then $I \setminus \{in_1, \dots, in_n\} = I$, hence the result.

Otherwise, we assume by induction that $\text{Redundant}(I \setminus \{in_1, \dots, in_n\}) \subseteq \text{Redundant}(I)$.

According to Lemma 4, redundancies can only decrease when performing a removal step. So, for every input $in_0 \in \text{Redundant}(I \setminus \{in_1, \dots, in_n, in_{n+1}\})$, we have:

$$\begin{aligned} 0 &< \text{redundancy}(in_0, I \setminus \{in_1, \dots, in_n, in_{n+1}\}) \\ &\leq \text{redundancy}(in_0, I \setminus \{in_1, \dots, in_n\}) \end{aligned}$$

So, $\text{Redundant}(I \setminus \{in_1, \dots, in_n, in_{n+1}\}) \subseteq \text{Redundant}(I \setminus \{in_1, \dots, in_n\}) \subseteq \text{Redundant}(I)$. \square

Since orders of removal steps contain inputs without repetition (Lemma 5), the index $\text{index}(in, \ell)$ of an input in in a removal order ℓ containing in is well defined. We leverage this to define subsequences.

Definition 1 (Subsequences). Let $[in_1, \dots, in_n]$ and $[in'_1, \dots, in'_m]$ be two orders of removal steps.

We say that $[in'_1, \dots, in'_m]$ is a *subsequence* of $[in_1, \dots, in_n]$ if $[in'_1, \dots, in'_m]$ can be obtained by removing some elements of $[in_1, \dots, in_n]$, i.e., if $\{in \in [in'_1, \dots, in'_m]\} \subseteq \{in \in [in_1, \dots, in_n]\}$ and, for any two inputs $in_i, in_j \in [in'_1, \dots, in'_m]$, if $\text{index}(in_i, [in_1, \dots, in_n]) \leq \text{index}(in_j, [in_1, \dots, in_n])$, then $\text{index}(in_i, [in'_1, \dots, in'_m]) \leq \text{index}(in_j, [in'_1, \dots, in'_m])$, where $\text{index}(in, \ell)$ is the index of input in in ℓ .

Lemma 7. *Let I be an input set and $[in_1, \dots, in_n], [in'_1, \dots, in'_m]$ be two orders of removal steps in I . If $[in_1, \dots, in_n] \in \text{ValidOrders}(I)$ and $[in'_1, \dots, in'_m]$ is a subsequence of $[in_1, \dots, in_n]$, then $[in'_1, \dots, in'_m] \in \text{ValidOrders}(I)$.*

Proof. The proof is done by induction on m .

If $m = 0$ then $[in'_1, \dots, in'_m] = [] \in \text{ValidOrders}(I)$.

Otherwise, we consider $[in'_1, \dots, in'_m, in'_{m+1}]$ and we assume by induction that $[in'_1, \dots, in'_m] \in \text{ValidOrders}(I)$.

Because $[in'_1, \dots, in'_m, in'_{m+1}]$ is a subsequence of $[in_1, \dots, in_n]$, there exists an index i such that $in_i = in'_{m+1}$. By definition of valid removal steps (§III-C), we have $in_i \in \text{Redundant}(I \setminus \{in_1, \dots, in_{i-1}\})$.

Moreover, $[in'_1, \dots, in'_m, in'_{m+1}]$ is a subsequence of $[in_1, \dots, in_i]$, so we have $\{in \in [in'_1, \dots, in'_m]\} \subseteq \{in \in [in_1, \dots, in_{i-1}]\}$. Hence, according to Lemma 6 we have $\text{Redundant}(I \setminus \{in_1, \dots, in_{i-1}\}) \subseteq \text{Redundant}(I \setminus \{in'_1, \dots, in'_m\})$. Therefore, $in'_{m+1} = in_i \in \text{Redundant}(I \setminus \{in'_1, \dots, in'_m\})$.

$[in'_1, \dots, in'_m] \in \text{ValidOrders}(I)$ and $in'_{m+1} \in \text{Redundant}(I \setminus \{in'_1, \dots, in'_m\})$, so, according to our definition for valid orders of removal steps (§III-C), we have $[in'_1, \dots, in'_m, in'_{m+1}] \in \text{ValidOrders}(I)$. \square

Finally, we prove in Proposition 2 that inputs in a valid order of removal steps can be rearranged in any different order and the rearranged order of removal steps is also valid.

Lemma 8 (Transposition of a Valid Order). *Let I be an input set.*

If $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$ and $1 \leq i < j \leq n$,

then $[in_1, \dots, in_j, \dots, in_i, \dots, in_n] \in \text{ValidOrders}(I)$

where in_i and in_j were exchanged and the other inputs are left unchanged.

Proof. The proof is done in four steps.

1) $[in_1, \dots, in_{i-1}, in_j]$ is a subsequence of $[in_1, \dots, in_{i-1}, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$. So, according to Lemma 7, we have $[in_1, \dots, in_{i-1}, in_j] \in \text{ValidOrders}(I)$. Note that the first step even applies to the case $i = 1$.

2) If $j = i + 1$, then one can go directly to the third step. Otherwise, for each $i < k < j$, we denote:

$$\begin{aligned} I_{i+1} &\stackrel{\text{def}}{=} I \setminus \{in_1, \dots, in_{i-1}\} \\ I_{k+1} &\stackrel{\text{def}}{=} I_k \setminus \{in_k\} \end{aligned}$$

and, for the sake of conciseness, $I_k^i = I_k \setminus \{in_i\}$ and $I_k^j = I_k \setminus \{in_j\}$.

We now prove by induction on $i \leq k \leq j-1$ that:

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_k] \in \text{ValidOrders}(I)$$

For the initialization $k = i$, we have from the first step:

$$[in_1, \dots, in_{i-1}, in_j] \in \text{ValidOrders}(I)$$

We now assume by induction hypothesis on $i \leq k < j-1$ that:

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_k] \in \text{ValidOrders}(I)$$

We first prove, for each $bl \in \text{Cover}(in_{k+1})$, that:

$$|\text{Inputs}(bl) \cap I_{k+1}^j| > 1$$

According to Lemma 5, $in_1, \dots, in_i, \dots, in_j, \dots, in_n$ are distinct. We consider two cases.

a) We assume $bl \in \text{Cover}(in_j)$.

Because $I_j^i = I \setminus \{in_1, \dots, in_{i-1}, in_i, in_{i+1}, \dots, in_{j-1}\}$ and $I_{k+1}^j = I \setminus \{in_1, \dots, in_{i-1}, in_{i+1}, \dots, in_k, in_j\}$, we have $I_j^i \cup \{in_{k+1}\} \subseteq I_{k+1}^j \cup \{in_j\}$. Thus, for each $bl \in \text{Cover}(in_j) \cap \text{Cover}(in_{k+1})$:

$$|\text{Inputs}(bl) \cap I_j^i| + 1 \leq |\text{Inputs}(bl) \cap I_{k+1}^j| + 1$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$\text{redundancy}(in_j, I_j^i) > 0$$

So, for each $bl \in \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_j^i| > 1$$

Thus, for each $bl \in \text{Cover}(in_j) \cap \text{Cover}(in_{k+1})$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^j| > 1$$

b) We assume $bl \notin \text{Cover}(in_j)$.

For each $bl \in \text{Cover}(in_{k+1}) \setminus (\text{Cover}(in_i) \cup \text{Cover}(in_j))$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^i| = |\text{Inputs}(bl) \cap I_{k+1}^j|$$

Moreover, for each $bl \in \text{Cover}(in_i) \setminus \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^i| + 1 = |\text{Inputs}(bl) \cap I_{k+1}^j|$$

So, for each $bl \in \text{Cover}(in_{k+1}) \setminus \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^i| \leq |\text{Inputs}(bl) \cap I_{k+1}^j|$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$\text{redundancy}(in_{k+1}, I_{k+1}^i) > 0$$

So, for each $bl \in \text{Cover}(in_{k+1})$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^i| > 1$$

Thus, for each $bl \in \text{Cover}(in_{k+1}) \setminus \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_{k+1}^j| > 1$$

So, we proved by case, for each $bl \in \text{Cover}(in_{k+1})$, that:

$$|\text{Inputs}(bl) \cap I_{k+1}^j| > 1$$

Thus, we have $\text{redundancy}(in_{k+1}, I_{k+1}^j) > 0$.

Moreover, according to Lemma 5, $in_1, \dots, in_i, \dots, in_j, \dots, in_n$ are distinct, so $in_{k+1} \in I_{k+1}^j$. Hence, $in_{k+1} \in \text{Redundant}(I_{k+1}^j)$. By induction hypothesis, we have:

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_k] \in \text{ValidOrders}(I)$$

So, according to the definition of valid removal steps (§III-C):

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_k, in_{k+1}] \in \text{ValidOrders}(I)$$

which concludes the induction step.

Therefore, we proved by induction:

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_{j-1}] \in \text{ValidOrders}(I)$$

3) We now prove that, for each $bl \in \text{Cover}(in_i)$, we have:

$$|\text{Inputs}(bl) \cap I_j^j| > 1$$

where I_j^i and I_j^j are defined in step two. We consider two cases.

a) We assume $bl \in \text{Cover}(in_j)$. Because $I_j^i \cup \{in_i\} = I_j^j \cup \{in_j\}$, for each $bl \in \text{Cover}(in_i) \cap \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_j^i| + 1 = |\text{Inputs}(bl) \cap I_j^j| + 1$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$\text{redundancy}(in_j, I_j^i) > 0$$

So, for each $bl \in \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_j^i| > 1$$

Thus, for each $bl \in \text{Cover}(in_i) \cap \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_j^j| > 1$$

b) We assume $bl \notin \text{Cover}(in_j)$. In that case, for each $bl \in \text{Cover}(in_i) \setminus \text{Cover}(in_j)$, we have:

$$|\text{Inputs}(bl) \cap I_j^i| + 1 = |\text{Inputs}(bl) \cap I_j^j|$$

We consider two cases.

i) There exists $i+1 \leq k \leq j-1$ such that $bl \in \text{Cover}(in_k)$. In that case, let k_{\max} be the largest. So, we have:

$$(\text{Inputs}(bl) \cap I_j^i) \cup \{in_{k_{\max}}\} = \text{Inputs}(bl) \cap I_{k_{\max}}^i$$

Hence:

$$|\text{Inputs}(bl) \cap I_j^i| + 1 = |\text{Inputs}(bl) \cap I_{k_{\max}}^i|$$

Thus:

$$|\text{Inputs}(bl) \cap I_j^j| = |\text{Inputs}(bl) \cap I_{k_{\max}}^i|$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$\text{redundancy}(in_{k_{\max}}, I_{k_{\max}}^i) > 0$$

So, because $bl \in \text{Cover}(in_{k_{\max}})$, we have:

$$|\text{Inputs}(bl) \cap I_{k_{\max}}^i| > 1$$

Finally:

$$|\text{Inputs}(bl) \cap I_j^j| > 1$$

ii) Otherwise, for each $i + 1 \leq k \leq j - 1$, we have $bl \notin \text{Cover}(in_k)$. Note that this is also the case if $j = i + 1$. In that case, we have:

$$(\text{Inputs}(bl) \cap I_j^i) \cup \{in_i\} = \text{Inputs}(bl) \cap I_i$$

where $I_i = I \setminus \{in_1, \dots, in_{i-1}\}$. So:

$$|\text{Inputs}(bl) \cap I_j^i| + 1 = |\text{Inputs}(bl) \cap I_i|$$

Thus:

$$|\text{Inputs}(bl) \cap I_j^j| = |\text{Inputs}(bl) \cap I_i|$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$\text{redundancy}(in_i, I_i) > 0$$

So, because $bl \in \text{Cover}(in_i)$, we have:

$$|\text{Inputs}(bl) \cap I_i| > 1$$

Finally:

$$|\text{Inputs}(bl) \cap I_j^j| > 1$$

This concludes the proof by case for $bl \in \text{Cover}(in_i) \setminus \text{Cover}(in_j)$. Therefore, we proved by case that, for each $bl \in \text{Cover}(in_i)$, we have:

$$|\text{Inputs}(bl) \cap I_j^j| > 1$$

Thus, we have $\text{redundancy}(in_i, I_j^j) > 0$.

Moreover, according to Lemma 5, $in_1, \dots, in_i, \dots, in_j, \dots, in_n$ are distinct, so $in_i \in I_j^j$. Thus, $in_i \in \text{Redundant}(I_j^j)$.

Finally, we have from the second step (or the first one, if $j = i + 1$):

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_{j-1}] \in \text{ValidOrders}(I)$$

Therefore, according to the definition of valid orders of removal steps (§III-C):

$$[in_1, \dots, in_{i-1}, in_j, in_{i+1}, \dots, in_{j-1}, in_i] \in \text{ValidOrders}(I)$$

4) Finally, if $j = n$ then the proof is complete. Otherwise, we prove by induction on $j \leq k \leq n$ that:

$$[in_1, \dots, in_j, \dots, in_i, in_{j+1}, \dots, in_k] \in \text{ValidOrders}(I)$$

For the initialization $k = j$, we have from the third step:

$$[in_1, \dots, in_j, \dots, in_i] \in \text{ValidOrders}(I)$$

We now assume by induction hypothesis on $j \leq k < n$ that:

$$[in_1, \dots, in_j, \dots, in_i, in_{j+1}, \dots, in_k] \in \text{ValidOrders}(I)$$

Because $[in_1, \dots, in_i, \dots, in_j, \dots, in_n] \in \text{ValidOrders}(I)$, we have:

$$in_{k+1} \in \text{Redundant}(I \setminus \{in_1, \dots, in_i, \dots, in_j, \dots, in_k\})$$

Moreover, we have:

$$\begin{aligned} & \{in_1, \dots, in_i, \dots, in_j, \dots, in_k\} \\ &= \{in_1, \dots, in_j, \dots, in_i, \dots, in_k\} \end{aligned}$$

So:

$$in_{k+1} \in \text{Redundant}(I \setminus \{in_1, \dots, in_j, \dots, in_i, \dots, in_k\})$$

Thus, according to the definition of valid orders of removal steps (§III-C):

$$\begin{aligned} & [in_1, \dots, in_j, \dots, in_i, in_{j+1}, \dots, in_k, in_{k+1}] \\ & \in \text{ValidOrders}(I) \end{aligned}$$

which concludes the induction step.

Therefore, we proved by induction the lemma:

$$[in_1, \dots, in_j, \dots, in_i, in_{j+1}, \dots, in_n] \in \text{ValidOrders}(I)$$

□

The issue with our definition of the gain (§III-C) is that, to compute $\text{gain}(I)$ of an input set I , one has to try all the possible order of removal steps to determine which ones are valid. If I contains n inputs and we consider orders of $0 \leq k \leq n$ removal steps, there are $n \times \dots \times (n - k + 1) = \frac{n!}{(n-k)!}$ possibilities to investigate, which is usually large.

Thus, we present an optimization to reduce the cost of computing the gain. First, we prove in Proposition 2 that, while the order of inputs matters to determine if an order of removal steps is valid, it does not matter anymore once we know the order is valid. In other words, inputs in a valid order of removal steps may be rearranged in any different order, and the rearranged order of removal steps would be valid. Formally, if $[in_1, \dots, in_n]$ is a valid order of removal steps, then $[in_{\sigma(1)}, \dots, in_{\sigma(n)}]$ is also a valid order, when σ denotes a permutation of the n inputs, i.e., the same inputs but (potentially) in a different order:

Proposition 2 (Permutation of a Valid Order). *Let I be an input set. If $[in_1, \dots, in_n] \in \text{ValidOrders}(I)$ and σ is a permutation on n elements, then $[in_{\sigma(1)}, \dots, in_{\sigma(n)}] \in \text{ValidOrders}(I)$.*

Proof. Every permutation of a finite set can be expressed as the product of transpositions [1] (p.60). Let m be the number of such transpositions for σ . We prove the result by induction on m .

If $m = 0$, then σ is the identity and we have by hypothesis:

$$[in_{\sigma(1)}, \dots, in_{\sigma(n)}] \in \text{ValidOrders}(I)$$

If $\sigma = \tau_{m+1} \circ \tau_m \circ \dots \circ \tau_1$, then we denote $\sigma' = \tau_m \circ \dots \circ \tau_1$. By induction hypothesis, we have:

$$[in_{\sigma'(1)}, \dots, in_{\sigma'(n)}] \in \text{ValidOrders}(I)$$

Then, by applying Lemma 8 to the transposition $(ij) = \tau_{m+1}$, we have:

$$[in_{\tau_{m+1} \circ \sigma'(1)}, \dots, in_{\tau_{m+1} \circ \sigma'(n)}] \in \text{ValidOrders}(I)$$

Finally, because $\sigma = \tau_{m+1} \circ \sigma'$, we conclude the induction step. \square

Since inputs in a valid order of removal steps are always distinct (Lemma 5) and their order does not make a difference (Proposition 2), such orders can simply be seen as sets even if, in practice, for IMPRO and MOCCO, orders of removal steps are implemented as lists. More precisely, two valid orders of removal steps are equivalent if they correspond to the same set of inputs. We leverage this property to reduce the number of orders of removal steps to consider, e.g., since $[in_1, in_2]$ is equivalent to $[in_2, in_1]$, we only have to check that $[in_1, in_2]$ is valid to know whether $[in_2, in_1]$ is valid or not. This means that, when inputs are provided in a given order $[in_1, \dots, in_n]$, these tools need only to check for valid orders of removal steps in increasing order (hence without backtracking on previous inputs), i.e., $[in_{i_1}, \dots, in_{i_m}]$, where $1 \leq i_1 < \dots < i_m \leq n$. We call *canonical orders* such orders of removal steps. For instance, $[in_2, in_4, in_7]$ is in canonical order, but $[in_4, in_2, in_7]$ is not. Thus, to save time when checking the validity of orders of removal steps, we consider only canonical orders.

Theorem 1 (Canonical Order). *Let $I = \{in_1, \dots, in_n\}$ be an input set. There exists $[in_{i_1}, \dots, in_{i_m}] \in \text{ValidOrders}(I)$ such that $1 \leq i_1 < \dots < i_m \leq n$ and:*

$$\sum_{1 \leq j \leq m} \text{cost}(in_{i_j}) = \text{gain}(I)$$

Proof. Let $[in_{i'_1}, \dots, in_{i'_m}] \in \text{ValidOrders}(I)$ be a valid order of removal steps with a maximal cumulative cost i.e., according to our definition of the gain (§III-C):

$$\sum_{1 \leq j \leq m} \text{cost}(in_{i'_j}) = \text{gain}(I)$$

If $m = 0$, then $[\]$ satisfies the theorem for a gain = 0.

Otherwise, we assume $m > 0$. According to Lemma 5, $[in_{i'_1}, \dots, in_{i'_m}]$ contains m distinct inputs. Let $\sigma \in S_m$ be the permutation such that:

$$[in_{\sigma(i'_1)}, \dots, in_{\sigma(i'_m)}] = [in_{i_1}, \dots, in_{i_m}]$$

with $1 \leq i_1 < \dots < i_m \leq n$.

According to Proposition 2, we have $[in_{i_1}, \dots, in_{i_m}] \in \text{ValidOrders}(I)$.

Moreover, because a permutation of the elements of a sum does not change the value of the sum, we have:

$$\sum_{1 \leq j \leq m} \text{cost}(in_{i_j}) = \sum_{1 \leq j \leq m} \text{cost}(in_{i'_j}) = \text{gain}(I)$$

Hence the theorem. \square

Theorem 1 allows us to focus on inputs in increasing order, instead of exploring all possible input orders. Thus, this optimization saves computation steps and makes computing the gain more tractable. More precisely, since there are $!k$ ways of ordering k selected inputs, there are “only” $\frac{!n}{!k!(n-k)}$ remaining possibilities to investigate, instead of $\frac{!n}{!(n-k)}$.

D. Local Dominance

In this section we prove, as presented in §VII-E, that the local dominance relation is local (Proposition 3), hence is faithful to its name, and that non locally-dominated inputs, as a whole, locally-dominate all the locally dominated inputs (Theorem 2). We start by the locality property.

Proposition 3 (Local Dominance is Local). *Let $in_1 \in I_{\text{search}}$ be a remaining input and $S \subseteq I_{\text{search}}$ be a subset of the remaining inputs. If $in_1 \sqsubseteq S$ then $in_1 \sqsubseteq S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\}$.*

Proof. We assume $in_1 \sqsubseteq S$. So, by definition of local dominance (§VII-E), we have $in_1 \notin S$, $\text{Cover}(in_1) \subseteq \text{Cover}(S)$, and $\text{cost}(in_1) \geq \text{cost}(S)$. First, because $in_1 \notin S$, we have $in_1 \notin S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\}$.

Moreover, for every $in_2 \in I_{\text{search}}$, if there is no overlap between in_1 and in_2 then, according to the overlapping relation (§VII-E), we have $\text{Cover}(in_1) \cap \text{Cover}(in_2) = \emptyset$. Thus:

$$\text{Cover}(in_1) \cap \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid \neg in_1 \sqcap in_2\}) = \emptyset$$

Hence, because $\text{Cover}(in_1) \subseteq \text{Cover}(S)$, we have:

$$\begin{aligned} & \text{Cover}(in_1) \\ &= \text{Cover}(in_1) \cap \text{Cover}(S) \\ &= \text{Cover}(in_1) \cap (\text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\}) \\ & \quad \cup \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid \neg in_1 \sqcap in_2\})) \\ &= (\text{Cover}(in_1) \cap \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\})) \\ & \quad \cup (\text{Cover}(in_1) \cap \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid \neg in_1 \sqcap in_2\})) \\ &= (\text{Cover}(in_1) \cap \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\})) \end{aligned}$$

So, $\text{Cover}(in_1) \subseteq \text{Cover}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\})$.

Finally, $\text{cost}(S) \geq \text{cost}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\})$. Thus, because $\text{cost}(in_1) \geq \text{cost}(S)$, we have $\text{cost}(in_1) \geq \text{cost}(S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\})$.

Therefore, $in_1 \sqsubseteq S \cap \{in_2 \in I_{\text{search}} \mid in_1 \sqcap in_2\}$. \square

We now prove that the local dominance relation is asymmetric (Proposition 5). First, because an input always covers some objectives, it can be locally dominated only by a non-empty input subset.

Lemma 9 (Non-Empty Local Dominance). *Let $in \in I_{\text{search}}$ be an input and $S \subseteq I_{\text{search}}$ be a subset of the remaining inputs. If $in \sqsubseteq S$ then $S \neq \emptyset$.*

Proof. The proof is done by contradiction. If $S = \emptyset$ then $\text{Cover}(S) = \emptyset$. Because $in \sqsubseteq S$ we have $\text{Cover}(in) \subseteq \text{Cover}(S)$. So, $\text{Cover}(in) = \emptyset$, which contradicts Lemma 1. \square

Second, when two inputs locally dominates each other, they are equivalent in the sense of the equivalence relation from §VII-D (two inputs are equivalent if they have the same coverage and cost), that we denote \equiv in the following. But, because we removed duplicates in §VII-D, this case can happen only if the inputs are the same, which is excluded by the definition of local dominance (§VII-E).

Lemma 10 (Local Dominance for Singletons is Asymmetric). *Let $in_1, in_2 \in I_{search}$ be two inputs. If $in_1 \sqsubseteq \{in_2\}$, then $in_2 \not\sqsubseteq \{in_1\}$.*

Proof. We assume by contradiction that $in_1 \sqsubseteq \{in_2\}$ and $in_2 \sqsubseteq \{in_1\}$.

Hence, we have: 1) $Cover(in_1) \subseteq Cover(in_2)$ and $Cover(in_2) \subseteq Cover(in_1)$, so $Cover(in_1) = Cover(in_2)$. 2) $cost(in_1) \geq cost(in_2)$ and $cost(in_2) \geq cost(in_1)$, so $cost(in_1) = cost(in_2)$. So, we have $in_1 \equiv in_2$.

Because we removed duplicates in §VII-D, we have $in_1 = in_2$. So, $in_1 \sqsubseteq \{in_1\}$, which contradicts in the definition of local dominance (§VII-E) that an input cannot locally dominates itself. \square

Note that this result is also useful to prove transitivity (Proposition 7).

Third, because the cost of an input is positive, if an input is locally dominated by several inputs, then they have a strictly smaller cost.

Lemma 11 (Cost Hierarchy). *Let $in_1 \in I_{search}$ be an input and $S \subseteq I_{search}$ be a subset of the remaining inputs. If $in_1 \sqsubseteq S$ and $|S| \geq 2$ then for every $in_2 \in S$ we have $cost(in_2) < cost(in_1)$.*

Proof. Because $in_1 \sqsubseteq S$ we have $cost(in_1) \geq cost(S)$. So, for every $in_2 \in S$ we have $cost(in_2) \leq cost(in_1)$. If $|S| \geq 2$, then the inequality is strict because, according to Lemma 2, the cost is positive $cost(in_2) > 0$. \square

Finally, we prove as expected that the local dominance relation is asymmetric.

Proposition 4 (Local Dominance for Subsets is Asymmetric). *Let $in_1, in_2 \in I_{search}$ be two inputs and $S_1, S_2 \subseteq I_{search}$ be two subsets of the remaining inputs. If $in_1 \sqsubseteq S_1$, $in_2 \in S_1$, and $in_2 \sqsubseteq S_2$, then $in_1 \not\sqsubseteq S_2$.*

Proof. The proof is made by contradiction. We assume $in_1 \in S_2$ and prove a contradiction in different cases for $|S_1|$ and $|S_2|$.

$|S_1| = 0$ or $|S_2| = 0$ are not possible, because this would contradict Lemma 9.

If $|S_1| = 1$ and $|S_2| = 1$, then $S_1 = \{in_2\}$ and $S_2 = \{in_1\}$. Thus, $in_1 \sqsubseteq \{in_2\}$ and $in_2 \sqsubseteq \{in_1\}$, which contradicts Lemma 10.

If $|S_1| = 1$ then $S_1 = \{in_2\}$. Because $in_1 \sqsubseteq S_1$, we have $cost(in_1) \geq cost(in_2)$. Moreover, because $in_2 \sqsubseteq S_2$ and $in_1 \in S_2$, if $|S_2| \geq 2$ then according to Lemma 11 we have $cost(in_1) < cost(in_2)$, hence the contradiction.

Because $in_1 \sqsubseteq S_1$ and $in_2 \in S_1$, if $|S_1| \geq 2$ then according to Lemma 11 we have $cost(in_2) < cost(in_1)$. Because $in_2 \sqsubseteq S_2$ and $in_1 \in S_2$, if $|S_2| \geq 2$ then according to Lemma 11 we have $cost(in_1) < cost(in_2)$. Hence the contradiction $cost(in_1) < cost(in_1)$. \square

Based on our definition of local dominance for subsets (§VII-E), we introduce the corresponding definition for inputs, in order to state more easily Corollary 1, then prove Theorem 2.

Definition 2 (Local Dominance). The input $in_1 \in I_{search}$ locally-dominates the input $in_2 \in I_{search}$, denoted $in_1 \hookrightarrow in_2$, if there exists a subset $S \subseteq I_{search}$ such that $in_1 \in S$ and $in_2 \sqsubseteq S$.

Proposition 5 (Local Dominance for Inputs is Asymmetric). *The \hookrightarrow relation is asymmetric i.e., for every $in_1, in_2 \in I_{search}$, if $in_2 \hookrightarrow in_1$, then $in_1 \not\hookrightarrow in_2$.*

Proof. If $in_2 \hookrightarrow in_1$ then there exists $S_1 \subseteq I_{search}$ such that $in_2 \in S_1$ and $in_1 \sqsubseteq S_1$. The proof is done by contradiction, assuming $in_1 \hookrightarrow in_2$. Hence, there exists $S_2 \subseteq I_{search}$ such that $in_1 \in S_2$ and $in_2 \sqsubseteq S_2$. According to Proposition 4 we have $in_1 \not\sqsubseteq S_2$, hence the contradiction. \square

We now prove that the local dominance relation is transitive (Proposition 7).

Proposition 6 (Local Dominance for Subsets is Transitive). *Let $in_1, in_2 \in I_{search}$ be two inputs and $S_1, S_2 \subseteq I_{search}$ be two subsets of the remaining inputs.*

If $in_1 \sqsubseteq S_1$, $in_2 \in S_1$, and $in_2 \sqsubseteq S_2$, then $in_1 \sqsubseteq (S_1 \setminus \{in_2\}) \cup S_2$.

Proof. Because $Cover(in_1) \subseteq Cover(S_1)$, $in_2 \in S_1$, and $Cover(in_2) \subseteq Cover(S_2)$, we have:

$$\begin{aligned} Cover(in_1) &\subseteq Cover(S_1 \setminus \{in_2\}) \cup Cover(S_2) \\ &= Cover((S_1 \setminus \{in_2\}) \cup S_2) \end{aligned}$$

Moreover, because $cost(in_1) \geq cost(S_1)$, $in_2 \in S_1$, and $cost(in_2) \geq cost(S_2)$, we have:

$$\begin{aligned} cost(in_1) &\geq cost(S_1 \setminus \{in_2\}) + cost(S_2) \\ &\geq cost((S_1 \setminus \{in_2\}) \cup S_2) \end{aligned}$$

Finally, because $in_1 \sqsubseteq S_1$ we have $in_1 \notin S_1$. We prove $in_1 \notin (S_1 \setminus \{in_2\}) \cup S_2$ by contradiction. If $in_1 \in (S_1 \setminus \{in_2\}) \cup S_2$ then, because $in_1 \notin S_1$, we have $in_1 \in S_2$. In that case, we have

$$\begin{aligned} cost(in_1) &\geq cost(S_1 \setminus \{in_2\}) + cost(S_2) \\ &= cost(S_1 \setminus \{in_2\}) + cost(S_2 \setminus \{in_1\}) + cost(in_1) \\ &\geq cost(in_1) \end{aligned}$$

Hence, by subtracting $cost(in_1)$, we have:

$$cost(S_1 \setminus \{in_2\}) + cost(S_2 \setminus \{in_1\}) = 0$$

Thus, according to Lemma 2, we have $S_1 = \{in_2\}$ and $S_2 = \{in_1\}$. Therefore $in_1 \sqsubseteq \{in_2\}$ and $in_2 \sqsubseteq \{in_1\}$, which contradicts Lemma 10. \square

Proposition 7 (Local Dominance for Inputs is Transitive). *The \hookrightarrow relation is transitive i.e., for every $in_1, in_2, in_3 \in I_{search}$, if $in_3 \hookrightarrow in_2$ and $in_2 \hookrightarrow in_1$, then $in_3 \hookrightarrow in_1$.*

Proof. If $in_3 \hookrightarrow in_2$ and $in_2 \hookrightarrow in_1$, then there exists $S_1, S_2 \subseteq I_{search}$ such that $in_2 \in S_1$, $in_1 \sqsubseteq S_1$, $in_3 \in S_2$, and $in_2 \sqsubseteq S_2$. According to Proposition 6, we have $in_1 \sqsubseteq (S_1 \setminus \{in_2\}) \cup S_2$. Hence, because $in_3 \in S_2 \subseteq (S_1 \setminus \{in_2\}) \cup S_2$ we have $in_3 \hookrightarrow in_1$. \square

We finally prove that the local dominance relation is acyclic.

Corollary 1 (Local Dominance is Acyclic). *There exists no cycle $in_1 \hookrightarrow \dots \hookrightarrow in_n \hookrightarrow in_1$ such that $in_1, \dots, in_n \in I_{search}$.*

Proof. The proof is done by contradiction, by assuming the existence of such a cycle $in_1 \hookrightarrow \dots \hookrightarrow in_n \hookrightarrow in_1$. According to Proposition 7, we have by transitivity that $in_1 \hookrightarrow in_1$. Then, according to Proposition 5, we have by asymmetry $in_1 \not\hookrightarrow in_1$. Hence the contradiction. \square

Finally, we use the transitivity (Proposition 7 and Proposition 6) and the acyclicity (Corollary 1) of the local-dominance relation to prove that non locally-dominated inputs, as a whole, locally-dominate all the locally dominated inputs.

Theorem 2 (Local Dominance Hierarchy). *For every locally-dominated input $in \in I_{search}$, there exists a subset $S \subseteq I_{search}$ of not locally-dominated inputs such that $in \sqsubseteq S$.*

Proof. For each input $in_0 \in I_{search}$ we denote:

$$LocDoms(in_0) \stackrel{\text{def}}{=} \{in \in I_{search} \mid in \hookrightarrow in_0\}$$

the input set which locally dominate in_0 .

Let $in_0 \in I_{search}$ be an input. We prove by induction on $LocDoms(in_0)$ that either in_0 is not locally dominated or there exists a subset $S \subseteq I_{search}$ of not locally-dominated inputs such that $in_0 \sqsubseteq S$.

For the initialization, we consider $LocDoms(in_0) = \emptyset$. In that case, in_0 is not locally dominated.

For the induction step, we consider $LocDoms(in_0) \neq \emptyset$, so in_0 is a locally dominated input. Hence, there exists $S_0 \subseteq I_{search}$ such that $in_0 \sqsubseteq S_0$. We consider two cases.

Either every input in S_0 is not locally dominated. In that case $S = S_0$ satisfies the inductive property.

Or there exists $n \geq 1$ input in_1, \dots, in_n in S_0 which are locally dominated. The rest of the proof is done in two steps.

First, let $1 \leq i \leq n$. We prove that $LocDoms(in_i) \subset LocDoms(in_0)$, with a strict inclusion.

Because $in_i \in S_0$ and $in_0 \sqsubseteq S_0$, we have $in_i \hookrightarrow in_0$. Hence, $in_i \in LocDoms(in_0)$.

Let $in \in LocDoms(in_i)$, so we have $in \hookrightarrow in_i$. So, according to Proposition 7, we have by transitivity $in \hookrightarrow in_0$, thus $in \in LocDoms(in_0)$. Hence, $LocDoms(in_i) \subseteq LocDoms(in_0)$.

According to Corollary 1, there is no cycle so $in_i \not\hookrightarrow in_i$. Hence, $in_i \notin LocDoms(in_i)$.

Finally, we have $LocDoms(in_i) \subset LocDoms(in_0)$ and $in_i \in LocDoms(in_0) \setminus LocDoms(in_i)$. Therefore $LocDoms(in_i) \subset LocDoms(in_0)$.

Because $LocDoms(in_i) \subset LocDoms(in_0)$, we can apply the induction hypothesis on in_i , which is locally dominated. Therefore, for each $1 \leq i \leq n$, there exists a subset $S_i \subseteq I_{search}$ of not locally-dominated inputs such that $in_i \sqsubseteq S_i$.

Second, we use these S_i to define by induction on $0 \leq m < n$ the following input sets:

$$\begin{aligned} I_0 &= S_0 \\ I_{m+1} &= (I_m \setminus \{in_{m+1}\}) \cup S_{m+1} \end{aligned}$$

and we prove by induction on $0 \leq m \leq n$ that $in_0 \sqsubseteq I_m$ and that either, $m < n$ and in_{m+1}, \dots, in_n are the only locally

dominated inputs in I_m , or $m = n$ and I_m contains no locally dominated input.

For the initialization, we have $I_0 = S_0$ and the inductive property is satisfied because $in_0 \sqsubseteq S_0$ and in_1, \dots, in_n are the only locally dominated inputs in S_0 .

For the induction step $m+1 \leq n$, we assume that $in_0 \sqsubseteq I_m$ and that in_{m+1}, \dots, in_n are the only locally dominated inputs in I_m .

Because $I_{m+1} = (I_m \setminus \{in_{m+1}\}) \cup S_{m+1}$ and S_{m+1} contains no locally-dominated input, we have that either, $m+1 < n$ and in_{m+2}, \dots, in_n are the only locally dominated inputs in I_{m+1} , or $m+1 = n$ and I_{m+1} contains no locally-dominated input.

Moreover, because $in_0 \sqsubseteq I_m$, $in_{m+1} \in I_m$, and $in_{m+1} \sqsubseteq S_{m+1}$, then according to Proposition 6, we have $in_0 \sqsubseteq (I_m \setminus \{in_{m+1}\}) \cup S_{m+1}$. Hence $in_0 \sqsubseteq I_{m+1}$, which concludes the induction step.

Therefore, $in_0 \sqsubseteq I_n$ and I_n contains no locally-dominated input. Thus, $S = I_n$ satisfies the inductive property on $LocDoms(in_0)$. Hence the claim for not locally-dominated inputs. \square

E. Dividing the Problem

In this section, we prove that redundancy updates are local (Lemma 12), that reductions on different connected components can be performed independently (Proposition 8), and finally that the gain can be independently computed on each connected component (Theorem 3). The main purpose of the section is to justify we can divide our problem into subproblems (§VII-F). We start by the locality.

Lemma 12 (Redundancy Updates are Local). *Let I be an input set and let $in_1, in_2 \in I$.*

If $redundancy(in_2, I \setminus \{in_1\}) \neq redundancy(in_2, I)$, then $in_1 \sqcap in_2$.

Proof. The proof is done by contraposition. If in_1 and in_2 do not overlap, then $Cover(in_1) \cap Cover(in_2) = \emptyset$. So, for every $bl \in Cover(in_2)$, we have $in_1 \notin Inputs(bl)$, and thus $Inputs(bl) \cap (I \setminus \{in_1\}) = Inputs(bl) \cap I$. Therefore, $redundancy(in_2, I \setminus \{in_1\}) = redundancy(in_2, I)$. \square

Then, we prove the independence of removal steps performed on different components.

Lemma 13 (Independent Redundancies). *Let I be an input set. For each connected component $C \in Comps(I)$, for each input $in_0 \in C$, and for each order of removal steps $[in_1, \dots, in_n]$, if $in_1, \dots, in_n \notin C$, then $redundancy(in_0, I \setminus \{in_1, \dots, in_n\}) = redundancy(in_0, I)$.*

Proof. $Comps(I)$ are the connected components for the redundant inputs in I , hence they are disjoint. Therefore, for each $in_0 \in C$, because $in_1, \dots, in_n \notin C$, we have that in_0 do not overlap with any of the in_1, \dots, in_n . Therefore, according to Lemma 12, we have $redundancy(in_0, I) = redundancy(in_0, I \setminus \{in_1\}) = \dots = redundancy(in_0, I \setminus \{in_1, \dots, in_n\})$. \square

Proposition 8 (Independent Removal Steps).

Let I be an input set and let $C_1, \dots, C_c \in \text{Comps}(I)$ denote c connected components. If, for each $0 \leq i \leq c$, there exists inputs $in_1^i, \dots, in_{n_i}^i \in C_i$ such that $[in_1^i, \dots, in_{n_i}^i] \in \text{ValidOrders}(I)$, then:

$$[in_1^1, \dots, in_{n_1}^1] + \dots + [in_1^c, \dots, in_{n_c}^c] \in \text{ValidOrders}(I)$$

where $+$ denotes sequence concatenation.

Proof. The proof is done by induction on c .

If $c = 0$, then $[in_1^1, \dots, in_{n_1}^1] + \dots + [in_1^c, \dots, in_{n_c}^c] = [] \in \text{ValidOrders}(I)$.

We assume by induction that $[in_1^1, \dots, in_{n_1}^1] + \dots + [in_1^c, \dots, in_{n_c}^c] \in \text{ValidOrders}(I)$.

We denote ℓ this order of removal steps and for the sake of simplicity we also denote $S = \{in_1^1, \dots, in_{n_1}^1, \dots, in_1^c, \dots, in_{n_c}^c\}$.

By induction, let $C_{c+1} \in \text{Comps}(I)$ be another connected component and let $in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1} \in C_{c+1}$ such that $[in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1}] \in \text{ValidOrders}(I)$.

We now prove that $\ell + [in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1}] \in \text{ValidOrders}(I)$.

This is done by proving by induction on $0 \leq j \leq n_{c+1}$ that:

$$\ell + [in_1^{c+1}, \dots, in_j^{c+1}] \in \text{ValidOrders}(I)$$

and that, for each $in_0 \in C_{c+1} \setminus S_j$, we have:

$$\text{redundancy}(in_0, I \setminus (S \cup S_j)) = \text{redundancy}(in_0, I \setminus S_j)$$

where $S_j = \{in_1^{c+1}, \dots, in_j^{c+1}\}$.

If $n_{c+1} = 0$, then $\ell + [in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1}] = \ell \in \text{ValidOrders}(I)$. Moreover, because C_{c+1} is disjoint with C_1, \dots, C_c , according to Lemma 13, performing the ℓ removal steps does not change the redundancies of inputs in C_{c+1} i.e., for each $in_0 \in C_{c+1}$, we have $\text{redundancy}(in_0, I \setminus S) = \text{redundancy}(in_0, I)$.

We now assume by induction that $\ell + [in_1^{c+1}, \dots, in_j^{c+1}] \in \text{ValidOrders}(I)$ and for each $in_0 \in C_{c+1} \setminus S_j$, we have $\text{redundancy}(in_0, I \setminus (S \cup S_j)) = \text{redundancy}(in_0, I \setminus S_j)$.

Because $[in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1}] \in \text{ValidOrders}(I)$, we have $in_{j+1}^{c+1} \in \text{Redundant}(I \setminus S_j)$. Moreover, $in_{j+1}^{c+1} \in C_{c+1}$ hence $in_{j+1}^{c+1} \notin S$. Therefore, using the induction hypothesis with $in_0 = in_{j+1}^{c+1}$, we have $in_{j+1}^{c+1} \in \text{Redundant}(I \setminus (S \cup S_j))$. Therefore, according to our definition of valid removal steps (§III-C), $\ell + [in_1^{c+1}, \dots, in_{j+1}^{c+1}] \in \text{ValidOrders}(I)$.

We denote $S_{j+1} = S_j \cup \{in_{j+1}^{c+1}\}$. Let $in_0 \in C_{c+1} \setminus S_{j+1}$.

To complete the induction step, we prove that:

$$\begin{aligned} \text{redundancy}(in_0, I \setminus (S \cup S_{j+1})) \\ = \text{redundancy}(in_0, I \setminus S_{j+1}) \end{aligned}$$

We remind that, for each input set X , we have:

$$\begin{aligned} \text{redundancy}(in_0, X) \\ = \min\{|\text{Inputs}(bl) \cap X| \mid bl \in \text{Cover}(in_0)\} - 1 \end{aligned}$$

We denote as critical the objectives contributing to the redundancy:

$$\begin{aligned} \text{CritSubCls}(in_0, X) \\ \stackrel{\text{def}}{=} \arg \min\{|\text{Inputs}(bl) \cap X| \mid bl \in \text{Cover}(in_0)\} \end{aligned}$$

Because $in_0 \in C_{c+1}$, we know that in_0 does not overlap with inputs in S so, for each $bl \in \text{Cover}(in_0)$, we have $\text{Inputs}(bl) \cap (I \setminus S) = \text{Inputs}(bl) \cap I$. Thus, by removing inputs of S_j from both sides, we have $\text{Inputs}(bl) \cap (I \setminus (S \cup S_j)) = \text{Inputs}(bl) \cap (I \setminus S_j)$. Therefore:

$$\text{CritSubCls}(in_0, I \setminus (S \cup S_j)) = \text{CritSubCls}(in_0, I \setminus S_j)$$

We consider two cases:

1) Either there exists $bl \in \text{CritSubCls}(in_0, I \setminus S_j)$ such that $in_{j+1}^{c+1} \in \text{Inputs}(bl)$. In that case:

$$\begin{aligned} |\text{Inputs}(bl) \cap (I \setminus (S \cup S_{j+1}))| \\ = |\text{Inputs}(bl) \cap (I \setminus (S \cup S_j))| - 1 \end{aligned}$$

$$|\text{Inputs}(bl) \cap (I \setminus S_{j+1})| = |\text{Inputs}(bl) \cap (I \setminus S_j)| - 1$$

Because bl is critical and a redundancy can decrease at most by 1 after a reduction (Lemma 4) so the cardinalities for other objectives cannot decrease below the previous redundancy minus one, we have:

$$\begin{aligned} \text{redundancy}(in_0, I \setminus (S \cup S_{j+1})) \\ = \text{redundancy}(in_0, I \setminus (S \cup S_j)) - 1 \end{aligned}$$

$$\text{redundancy}(in_0, I \setminus S_{j+1}) = \text{redundancy}(in_0, I \setminus S_j) - 1$$

2) Or for each $bl \in \text{CritSubCls}(in_0, I \setminus S_j)$, we have $in_{j+1}^{c+1} \notin \text{Inputs}(bl)$. In that case:

$$\begin{aligned} |\text{Inputs}(bl) \cap (I \setminus (S \cup S_{j+1}))| \\ = |\text{Inputs}(bl) \cap (I \setminus (S \cup S_j))| \end{aligned}$$

$$|\text{Inputs}(bl) \cap (I \setminus S_{j+1})| = |\text{Inputs}(bl) \cap (I \setminus S_j)|$$

Because bl is critical and a redundancy can decrease at most by 1 after a reduction (Lemma 4) so the cardinalities for non-critical objectives cannot decrease below the previous redundancy we have:

$$\begin{aligned} \text{redundancy}(in_0, I \setminus (S \cup S_{j+1})) \\ = \text{redundancy}(in_0, I \setminus (S \cup S_j)) \end{aligned}$$

$$\text{redundancy}(in_0, I \setminus S_{j+1}) = \text{redundancy}(in_0, I \setminus S_j)$$

In both cases, by induction hypothesis $\text{redundancy}(in_0, I \setminus (S \cup S_j)) = \text{redundancy}(in_0, I \setminus S_j)$, hence we have $\text{redundancy}(in_0, I \setminus (S \cup S_{j+1})) = \text{redundancy}(in_0, I \setminus S_{j+1})$.

This completes the induction on $0 \leq j \leq n_{c+1}$. In particular, we proved that $\ell + [in_1^{c+1}, \dots, in_{n_{c+1}}^{c+1}] \in \text{ValidOrders}(I)$, which completes the induction on c , hence the result. \square

We now prove that the gain can be independently computed on each connected component (Theorem 3). To do so, we have to prove intermediate results, including Proposition 9.

Lemma 14 (Redundancy within a connected component). *Let I be an input set. For each connected component $C \in \text{Comps}(I)$, for each input $in_0 \in C$, and for each order of removal steps $[in_1, \dots, in_n]$, if $in_1, \dots, in_n \in C$, then:*

$$\begin{aligned} \text{redundancy}(in_0, C \setminus \{in_1, \dots, in_n\}) \\ = \text{redundancy}(in_0, I \setminus \{in_1, \dots, in_n\}) \end{aligned}$$

Proof. We remind that, for each input set X , we have:

$$\begin{aligned} & \text{redundancy}(in_0, X) \\ &= \min\{|Inputs(bl) \cap X| \mid bl \in Cover(in_0)\} - 1 \end{aligned}$$

For the sake of simplicity, we denote $X_C = C \setminus \{in_1, \dots, in_n\}$ and $X_I = I \setminus \{in_1, \dots, in_n\}$. Let $in_0 \in C$ and $bl \in Cover(in_0)$.

Because $C \in Comps(I)$, we have $C \subseteq I$, hence $X_C \subseteq X_I$. So $Inputs(bl) \cap X_C \subseteq Inputs(bl) \cap X_I$. Moreover, for each $in \in Inputs(bl) \cap X_I$ we have $bl \in Cover(in_0) \cap Cover(in)$, so $in_0 \sqcap in$, and thus $in \in X_C$.

Therefore, $Inputs(bl) \cap X_C = Inputs(bl) \cap X_I$.

Hence the result $\text{redundancy}(in_0, X_C) = \text{redundancy}(in_0, X_I)$. \square

Corollary 2 (Component Validity). *Let I be an input set. For each connected component $C \in Comps(I)$ and for each order of removal steps $[in_1, \dots, in_n]$, if $[in_1, \dots, in_n] \in ValidOrders(C)$, then $[in_1, \dots, in_n] \in ValidOrders(I)$.*

Proof. The proof is done by induction on n .

For the initialization $n = 0$, we have $[in_1, \dots, in_n] = [] \in ValidOrders(I)$.

For the induction step, we consider $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(C)$ and we assume by induction that $[in_1, \dots, in_n] \in ValidOrders(I)$.

Because $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(C)$, according to our definition of redundant inputs and removal steps (§III-C), we have $in_1, \dots, in_n \in C$. Hence, according to Lemma 14, we have for each $in_0 \in C$:

$$\begin{aligned} & \text{redundancy}(in_0, C \setminus \{in_1, \dots, in_n\}) \\ &= \text{redundancy}(in_0, I \setminus \{in_1, \dots, in_n\}) \end{aligned}$$

Thus, because $C \subseteq I$, according to the definition of redundancy (§III-C) we have:

$$\begin{aligned} & \text{Redundant}(C \setminus \{in_1, \dots, in_n\}) \\ & \subseteq \text{Redundant}(I \setminus \{in_1, \dots, in_n\}) \end{aligned}$$

Because $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(C)$, we have $in_{n+1} \in \text{Redundant}(C \setminus \{in_1, \dots, in_n\})$ by definition of removal steps (§III-C). Thus, $in_{n+1} \in \text{Redundant}(I \setminus \{in_1, \dots, in_n\})$.

By induction hypothesis $[in_1, \dots, in_n] \in ValidOrders(I)$.

Therefore $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(I)$, which completes the induction step. \square

Proposition 9 (Reductions withing a connected component). *Let I be an input set. For each connected component $C \in Comps(I)$ and for each order of removal steps $[in_1, \dots, in_n]$, if $[in_1, \dots, in_n] \in ValidOrders(I)$ and $in_1, \dots, in_n \in C$, then $[in_1, \dots, in_n] \in ValidOrders(C)$.*

Proof. The proof is done by induction on n .

If $n = 0$ then $[in_1, \dots, in_n] = [] \in ValidOrders(C)$.

Otherwise, we consider $[in_1, \dots, in_n, in_{n+1}]$ with $in_1, \dots, in_n, in_{n+1} \in C$ and we assume by induction that $[in_1, \dots, in_n] \in ValidOrders(C)$.

We assume $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(I)$, so $in_{n+1} \in \text{Redundant}(I \setminus \{in_1, \dots, in_n\})$.

According to Lemma 14 applied to inputs $in_0 \in C \setminus \{in_1, \dots, in_n\}$ with redundancy > 0 , we have:

$$\begin{aligned} & \text{Redundant}(C \setminus \{in_1, \dots, in_n\}) \\ &= \text{Redundant}(I \setminus \{in_1, \dots, in_n\}) \end{aligned}$$

Hence, $in_{n+1} \in \text{Redundant}(C \setminus \{in_1, \dots, in_n\})$.

Thus, because $[in_1, \dots, in_n] \in ValidOrders(C)$, according to our definition of valid removal steps (§III-C), we have the result $[in_1, \dots, in_n, in_{n+1}] \in ValidOrders(C)$. \square

Finally, we conclude the section by proving the claim on Theorem 3, which is used to divide our problem into subproblems (§VII-F).

Theorem 3 (Divide the Gain). *For each input set I :*

$$\text{gain}(I) = \sum_{C \in Comps(I)} \text{gain}(C)$$

Proof. Let $[in_1, \dots, in_n] \in ValidOrders(I)$ be a valid order of removal steps such that its cumulative cost $\sum_{1 \leq i \leq n} \text{cost}(in_i) = \text{gain}(I)$ is maximal (§III-C), and let $Comps(I) = \{C_1, \dots, C_c\}$ denote the connected components, without a particular order.

We denote $[in_1, \dots, in_n]_{C_j}$ the largest sublist (Definition 1) of $[in_1, \dots, in_n]$ containing only inputs in the connected component C_j . Because $[in_1, \dots, in_n] \in ValidOrders(I)$, according to Lemma 7 we have $[in_1, \dots, in_n]_{C_j} \in ValidOrders(I)$ as well. So, according to Proposition 8:

$$[in_1, \dots, in_n]_{C_1} + \dots + [in_1, \dots, in_n]_{C_c} \in ValidOrders(I)$$

Because the connected components form a partition of the redundant inputs, the cumulative cost of this order of removal steps is the same as $[in_1, \dots, in_n]$:

$$\sum_{1 \leq j \leq c} \sum_{in \in [in_1, \dots, in_n]_{C_j}} \text{cost}(in) = \sum_{1 \leq i \leq n} \text{cost}(in_i) = \text{gain}(I)$$

We now consider any connected component C_{j_1} and we prove that:

$$\sum_{in \in [in_1, \dots, in_n]_{C_{j_1}}} \text{cost}(in) = \text{gain}(C_{j_1})$$

The proof is done in two steps.

First, note that because $[in_1, \dots, in_n]_{C_{j_1}} \in ValidOrders(I)$ and contains only inputs in C_{j_1} , according to Proposition 9 we have $[in_1, \dots, in_n]_{C_{j_1}} \in ValidOrders(C_{j_1})$ as well.

Second, we prove by contradiction that $[in_1, \dots, in_n]_{C_{j_1}}$ has a maximal cumulative cost in C_{j_1} . We assume by contradiction that there exists a valid order of removal steps $[in'_1, \dots, in'_{n'}] \in ValidOrders(C_{j_1})$ such that:

$$\sum_{in \in [in'_1, \dots, in'_{n'}]} \text{cost}(in) > \sum_{in \in [in_1, \dots, in_n]_{C_{j_1}}} \text{cost}(in)$$

Because $[in'_1, \dots, in'_{n'}] \in ValidOrders(C_{j_1})$, according to Corollary 2 we have $[in'_1, \dots, in'_{n'}] \in ValidOrders(I)$ as well. So, according to Proposition 8, we have:

$$[in_1, \dots, in_n]_{C_1} + \dots + [in_1, \dots, in_n]_{C_{j_1-1}}$$

$$+ [in'_1, \dots, in'_{n'}] + [in_1, \dots, in_n]_{C_{j_1+1}} + \dots \\ + [in_1, \dots, in_n]_{C_c} \in \text{ValidOrders}(I)$$

The cumulative cost of this valid order of removal steps is thus larger than the cumulative cost of $[in_1, \dots, in_n]_{C_1} + \dots + [in_1, \dots, in_n]_{C_c}$:

$$\sum_{in \in [in'_1, \dots, in'_{n'}]} \text{cost}(in) + \sum_{1 \leq j \leq c \wedge j \neq j_1} \sum_{in \in [in_1, \dots, in_n]_{C_j}} \text{cost}(in) \\ > \sum_{1 \leq j \leq c} \sum_{in \in [in_1, \dots, in_n]_{C_j}} \text{cost}(in) = \text{gain}(I)$$

which contradicts the maximality of $\text{gain}(I)$.

Hence, $[in_1, \dots, in_n]_{C_{j_1}} \in \text{ValidOrders}(C_{j_1})$ has a maximal cumulative cost in C_{j_1} . So, according to the definition of the gain (§III-C), we have for any connected component C_{j_1} :

$$\sum_{in \in [in_1, \dots, in_n]_{C_{j_1}}} \text{cost}(in) = \text{gain}(C_{j_1})$$

Therefore, we have the result:

$$\text{gain}(I) = \sum_{1 \leq j \leq c} \sum_{in \in [in_1, \dots, in_n]_{C_j}} \text{cost}(in) = \sum_{1 \leq j \leq c} \text{gain}(C_j)$$

□

F. Genetic Search

In this section, we prove desirable properties satisfied for each generation by roofers (Theorem 4) and misers (Theorem 5) during the genetic search (Section VIII). We start by roofers.

Theorem 4 (Invariant of the Roofers). *For every generation n , we have:*

$$|\text{Roofers}(n)| = n_{\text{size}} \\ \min\{\text{cost}(I) \mid I \in \text{Roofers}(n)\} \\ = \min\{\text{cost}(I) \mid I \in \bigcup_{0 \leq m \leq n} \text{Roofers}(m)\}$$

and for every $I \in \text{Roofers}(n)$, we have:

- I is reduced (in the sense of §III-C)
- $\text{Cover}(I) = \text{Coverage}_{\text{obj}}(C)$

Proof. There are n_{size} individuals in the initial roofer population (§VIII-B). Moreover, the procedure detailed above to update the roofer population ensures that an offspring can only take the place of an existing roofer. Hence, the number of roofers does not change over generations.

In the initial population, a roofer I is always replaced by its reduced counterpart $\text{reduce}(I)$ after adding an input (§VIII-B). Moreover, after mutation (§VIII-E) each offspring I is replaced by its reduced counterpart $\text{reduce}(I)$ before determining if it is accepted in the roofer population or rejected. Hence, for each generation, each roofer is reduced.

Roofers in the initial population are built so that they cover all the objectives (§VIII-B). Moreover, in the above procedure, an offspring can be added to the roofer population only if it covers all the objectives. Hence, for each generation, each roofer covers all the objectives.

Finally, in the above procedure one can remove an individual from the roofer population only if a less or equally costly roofer is found. Hence, the minimal cost amongst the roofers can only remain the same or decrease over generations. Therefore, the minimal cost in the last generation is the minimal cost encountered so far during the search. □

To ease the proof for misers (Theorem 5), we first prove in Lemma 16 that a miser I_0 can be removed between generation n and generation $n+1$ only by a dominating miser I . Then, we prove in Corollary 3 that being dominated is carried from generation to generation.

Lemma 15 (Transitivity of Pareto Dominance). *The \succ relation (§II-E) is transitive i.e., for every input sets I_1, I_2, I_3 , if $I_1 \succ I_2$ and $I_2 \succ I_3$, then $I_1 \succ I_3$.*

Proof. For each $0 \leq i \leq n$, we have $f_i(I_1) \leq f_i(I_2)$ and $f_i(I_2) \leq f_i(I_3)$, so $f_i(I_1) \leq f_i(I_3)$. Moreover, there exists $0 \leq i_1 \leq n$ such that $f_{i_1}(I_1) < f_{i_1}(I_2) \leq f_{i_1}(I_3)$ and there exists $0 \leq i_2 \leq n$ such that $f_{i_2}(I_2) < f_{i_2}(I_3)$. i_1 and i_2 can be the same or distinct. In any case, we have $f_{i_1}(I_1) < f_{i_1}(I_3)$ and $f_{i_2}(I_1) < f_{i_2}(I_3)$. □

Lemma 16. *For every generation n , if $I_0 \in \text{Misers}(n) \setminus \text{Misers}(n+1)$, then there exists $I \in \text{Misers}(n+1)$ such that $I \succ I_0$, where \succ is the domination relation from (§II-E).*

Proof. In the miser population update (§VIII-F), I_0 can be removed from the miser population only if there exists a miser candidate I_1 such that $I_1 \succ I_0$. If I_1 is accepted in the population then $I = I_1$ satisfies the lemma.

Otherwise, I_1 is rejected only because there exists a miser I_2 such that $I_2 \succ I_1$. Hence, according to Lemma 15, we have by transitivity that $I_2 \succ I_0$. Note that there is at most two miser candidates per generation. If either I_1 is the only miser candidate or there exists a second miser candidate I_3 which does not dominate I_2 , then I_2 is present in the next generation and $I = I_2$ satisfies the lemma.

Otherwise, there exists a second candidate I_3 such that $I_3 \succ I_2$. Hence, according to Lemma 15, we have by transitivity that $I_3 \succ I_0$. If I_3 is accepted in the population then $I = I_3$ satisfies the lemma.

Otherwise, I_3 is rejected only because there exists a miser I_4 such that $I_4 \succ I_3$. Hence, according to Lemma 15, we have by transitivity that $I_4 \succ I_0$. Because there is at most two miser candidates per generation, I_4 cannot be removed by another candidate. Therefore, I_4 is present in the next generation and $I = I_4$ satisfies the lemma. □

Corollary 3 (Dominance across Generations). *For every generation n and for every $I_1 \in \text{Misers}(n)$, if there exists $0 \leq m \leq n$ and $I_2 \in \text{Misers}(m)$ such that $I_2 \succ I_1$, then there exists $I_3 \in \text{Misers}(n)$ such that $I_3 \succ I_1$.*

Proof. The proof is done by induction on $n - m$.

If $I_2 \in \text{Misers}(n)$ then $I_2 = I_3$ satisfies the corollary.

Otherwise, there exists a generation $m < g \leq n$ where I_2 was removed. In that case, according to Lemma 16, there exists $I_3 \in \text{Misers}(g)$ such that $I_3 \succ I_2$. Hence, according to Lemma 15, we have by transitivity that $I_3 \succ I_1$. Finally,

because $n - g < n - m$, we have the result using the induction hypothesis on I_3 . \square

We finally prove, as expected, the properties satisfied by misers on each generation.

Theorem 5 (Invariant of the Misers). *For every generation n , for every $I_1 \in \text{Misers}(n)$, we have:*

- I_1 is reduced (in the sense of §III-C)
- $\text{Cover}(I) \subset \text{Coverage}_{obj}(C)$ (the inclusion is strict)
- There exists no $I_2 \in \bigcup_{0 \leq m \leq n} \text{Misers}(m)$ such that $I_2 \succ I_1$.

Proof. The initial miser population is empty (§VIII-B), hence the properties trivially hold.

After mutation (§VIII-E) each offspring I is replaced by its reduced counterpart $\text{reduce}(I)$ before determining if it is accepted in the miser population or rejected. Hence, for each generation, each miser is reduced.

In the miser population update (§VIII-F), an offspring can be a candidate to the miser population only if it does not cover all the objectives.

Finally, the last property is proved for a miser $I_1 \in \text{Misers}(n)$. Let n' be the first generation when I_1 was accepted, so we have $n' \leq n$ and $I_1 \in \text{Misers}(n')$.

The proof is done by contradiction, assuming that there existed a generation $0 \leq m \leq n$ and a miser $I_2 \in \text{Misers}(m)$ such that $I_2 \succ I_1$. Let m' be the first generation when I_2 was accepted, so we have $m' \leq m$ and $I_2 \in \text{Misers}(m')$. We consider two cases.

If $m' \leq n'$ then, according to Corollary 3, there exists $I_3 \in \text{Misers}(n')$ such that $I_3 \succ I_1$. This contradicts the above procedure, because if I_1 was dominated it would not have been accepted in $\text{Misers}(n')$.

If $m' > n'$ then, because $I_2 \succ I_1$, according to the above procedure I_1 is removed so $I_1 \notin \text{Misers}(m')$. But $m' \leq m \leq n$ and $I_1 \in \text{Misers}(n)$, so I_1 was added between m' and n . Let g be the first generation when this occurred.

In that case we have $I_1 \in \text{Misers}(g)$, $0 \leq m' \leq g$, $I_2 \in \text{Misers}(m')$, and $I_2 \succ I_1$. So, according to Corollary 3, there exists $I_3 \in \text{Misers}(g)$ such that $I_3 \succ I_1$, which contradicts the fact that I_1 was accepted in $\text{Misers}(g)$. \square

REFERENCES

- [1] M. J. Hall, *The Theory of Groups*. USA: MacMillan, 1959.