



# Keeping Secrets: Multi-objective Genetic Improvement for Detecting and Reducing Information Leakage

Ibrahim Mesecan  
Iowa State University  
Ames, IA, USA  
imesecan@iastate.edu

Myra B. Cohen  
Iowa State University  
Ames, IA, USA  
mcohen@iastate.edu

Daniel Blackwell, David Clark  
University College London  
London, UK  
david.clark@ucl.ac.uk  
daniel.blackwell.14@ucl.ac.uk

Justyna Petke  
University College London  
London, UK  
j.petke@ucl.ac.uk

## ABSTRACT

Information leaks in software can unintentionally reveal private data, yet they are hard to detect and fix. Although several methods have been proposed to detect leakage, such as static verification-based approaches, they require specialist knowledge, and are time-consuming. Recently, we introduced HyperGI, a dynamic, hypertest-based approach that can detect and produce potential fixes for hyperproperty violations. In particular, we focused on violations of the noninterference property, as it results in information flow leakage. Our instantiation of HyperGI was able to detect and reduce leakage in three small programs. Its fitness function tried to balance information leakage and program correctness but, as we pointed out, there may be tradeoffs between keeping program semantics and reducing information leakage that require developer decisions.

In this work we ask if it is possible to automatically detect and repair information leakage in more realistic programs without requiring specialist knowledge. We instantiate a multi-objective version of HyperGI in a tool, called LeakReducer, which explicitly encodes the tradeoff between program correctness and information leakage. We apply LeakReducer to six leaky programs, including the well-known Heartbleed bug. LeakReducer is able to detect leakage in all, in contrast to state-of-the-art fuzzers, detecting leakage in only two programs. Moreover, LeakReducer is able to reduce leakage in all subjects, with comparable results to previous work, while scaling to much larger software.

## CCS CONCEPTS

• **Software and its engineering** → **Search-based software engineering**; **Search-based software engineering**; • **Security and privacy** → **Software security engineering**.



This work is licensed under a Creative Commons Attribution International 4.0 License.

ASE '22, October 10–14, 2022, Rochester, MI, USA  
© 2022 Copyright held by the owner/author(s).  
ACM ISBN 978-1-4503-9475-8/22/10.  
<https://doi.org/10.1145/3551349.3556947>

## KEYWORDS

Genetic Improvement, Information Leakage, Search-Based Software Engineering, Automated Program Repair

### ACM Reference Format:

Ibrahim Mesecan, Daniel Blackwell, David Clark, Myra B. Cohen, and Justyna Petke. 2022. Keeping Secrets: Multi-objective Genetic Improvement for Detecting and Reducing Information Leakage. In *37th IEEE/ACM International Conference on Automated Software Engineering (ASE '22)*, October 10–14, 2022, Rochester, MI, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3551349.3556947>

## 1 INTRODUCTION

Information leakage from programs has led to high profile security bugs such as the Heartbleed bug.<sup>1</sup> Typically, information leaks from a program either when it contains information flow control (IFC) errors or when a data structure such as a buffer or stack can be accessed in an unbounded manner. The Heartbleed bug was due to a problem in the OpenSSL cryptographic library. When pinged with a malformed query, it was possible to read past the buffer and return unencrypted data from the server's memory, hence providing a backdoor for eavesdropping on network traffic. This bug existed for several years in a library used by programs and servers worldwide. While it represents a common type of information leakage, more subtle IFC leaks can occur, and these may also lead to exposure of private information.

Take, for instance, the program below which accepts an integer variable `var` and returns 0, 1 or 2. It has two predicates, one which compares `var` against a variable called `magic_number` and one that uses the value of a secret variable (`protected_var`).

```
int leaking_secrets(int var){
    ...
    if(var > magic_number){
        leak_info=2;
    }
    else if(var < protected_var){
        leak_info=1;
    }
    else
        leak_info=0;
    ...
    return leak_info; }
```

<sup>1</sup><https://heartbleed.com/>

An important part of information leakage is setting a security policy (beyond the scope of this work). Based on a given security policy, we can assume that 1) `var` and `magic_number` are publicly known (a.k.a low security variables), and 2) `protected_var` is considered high security; it is a secret. Given this policy, it is not hard to see that repeated querying of this program with different values of `var` can expose information about `protected_var`. Let's assume `magic_number=6` and `protected_var=5`. If we run the program repeatedly using the inputs 8, 3, 5, we get different return values of 2, 1, 0. Using inputs 0, 1, 2, we get return values of 1, 1, 1. These results give us information about the content of `protected_var`, i.e. information is leaking. As we can see, the value of `magic_number` could impact the visibility of this information. Even in this simple program the ability to discover the leak dynamically depends on multiple variables (and control flows). If, for instance, `magic_number` is set to a large negative value, then for most inputs of `var` no information about `protected_var` will be revealed. The program always returns 2.

Finding information leakage in programs is non-trivial, with the most common approach being static analysis [24, 38]. There have also been some combined static/dynamic techniques [34], but these have not necessarily been applied to real-world programs nor can they provide patch suggestions if the problem is found. Moreover, as we show in our study, current techniques may be able to detect information leakage related to memory overflows, but they may still miss those related to a program's control flow as is exemplified in the example program `leaking_secrets`. Mechtav et al. [30] were able to use automated program repair to fix the Heartbleed bug suggesting automated ways to handle these types of faults. However, that work assumes the leakage would break program functionality and requires failing test cases. As in our earlier work [32], we call such tests *functional tests*. Information leakage can occur even if a program passes all functional tests, making information leakage difficult to detect and fix.

Previously, we presented a general approach, called HyperGI, which uses hypertests and genetic improvement to repair software's hyperproperty violations. We implemented an instantiation of HyperGI for the confidentiality problem: i.e. detection and repair of information leakage in programs [32]. While this work made advances in reducing information leakage it was preliminary. Our implementation was applied only to small functions (less than 40 lines of code) and results were obtained in a partially automated setting, since a user needed to insert some domain knowledge (i.e., extra variables with their types) to improve search. Nor did it allow any flexibility in trading semantic invariance for leak repair.

A key finding of our previous work was that the patches produced presented a tradeoff between preserving original functionality (as exposed by functional tests) and reducing information leakage (as exposed by hypertests). Returning to our example above, the only way to reduce information leakage completely, would be to change the predicate related to `protected_var`. However, that would change the initial, intended program and some functional tests are likely to fail. This suggests repairing information flow leakage should be viewed as a multi-objective problem. We should consider the option of balancing the leakage of secrets with the need for particular program behavior.

In this paper, we present an automated multi-objective framework called LeakReducer for estimating and reducing information leakage. It requires only a program and its security policy. Moreover, we improve upon our previous work [32] by: (1) using multiple functional test sets as input; (2) using automatically derived repair ingredients; (3) including both single and multi-objective search strategies; and (4) scaling to real programs. We have evaluated LeakReducer on the prior subjects from our previous work as well as on three other programs, including two modules from the OpenSSL library. We were able to both detect and reduce leakage in two files, each around 1,000 lines of code, including the original Heartbleed bug. One of our detected leaks in a real-world program turned out to be a false positive (confirmed by developers), but it points to a different use case and class of information leakage faults which we aim to study as future work. Furthermore, we examine the quality of Pareto fronts for different multi-objective algorithms, and explore a few interesting patches in depth.

In summary, the contributions of this work are:

- (1) LeakReducer, a multi-objective instantiation and implementation of HyperGI for finding and reducing information leakage;
- (2) An empirical study demonstrating the effectiveness of LeakReducer; and
- (3) An evaluation of patch diversity produced by LeakReducer, detailing tradeoffs between change in original semantics of the given program and reduction in information leakage.

The rest of this paper is structured as follows: We present background and related work in the next section. We then present LeakReducer in Section 3. We follow this with our study and results (Sections 4-5). We conclude and suggest future work in Section 6.

## 2 BACKGROUND AND RELATED WORK

The information flow control problem of maintaining data confidentiality across program executions has a long history, part of which was surveyed in the early 2000s by Sabelfeld and Myers [37]. Software should be designed so that it obeys a security flow policy via a *noninterference* property, namely lower security users should not be aware of the actions of higher security users [21]. The security flow policies we consider are expressed using Lattice Based Access Control (LBAC) [18] where a subject can only access an object if the security level of the subject is greater than or equal to the object. Security levels are expressed in a lattice, a partial order where least upper bounds and greatest lower bounds are defined for any set. Such models are consistent with the “no write down, no read up” principles of Bell-LaPadula access control [9]. In this paper, all security policies are expressed in terms of a two point, High-Low labels lattice. Although LBAC allows complex multi-level policies, this lattice is sufficient to express the security policies for our programs under experimentation. Subjects and objects are mapped to the two points in this lattice and the noninterference property can be formulated purely in terms of these lattice points as follows.

We partition data containers, e.g., a memory address or a variable in an imperative program, and users using High/Low labels, then use the partition to formally define the noninterference property for the program and security policy pair. We say that a pair of states,  $s_1, s_2$ , are Low equivalent,  $s_1 \equiv_L s_2$ , if the values in the

data containers labelled Low are the same. Then we can define the property:

**DEFINITION 1 (NONINTERFERENCE PROPERTY).** *A program  $P$  satisfies the noninterference property for the High-Low security policy if for every pair of initial states  $s_1, s_2$ ,  $s_1 \equiv_L s_2 \Rightarrow P(s_1) \equiv_L P(s_2)$ .*

In this flow and termination insensitive definition, noninterference means the program will map all Low equivalent initial state pairs to Low equivalent final state pairs.

Noninterference is not a property of single executions but of *pairs* of executions. Clarkson and Schneider generalized this idea, calling such program properties *hyperproperties* as they are only expressible using sets of sets of executions [14]. The “preservation of Low equivalence” property partitions all distinct pairs of executions that begin in Low equivalent states into the set of pairs that preserve the Low equivalence in the final states and the set of pairs that do not. The noninterference property effectively says that the latter set is empty. If it is not, information leaks from the High labelled parts of initial states to the Low labelled parts of final states. This last observation leads to what Kinder called *Hypertests* [25], using pairs of Low equivalent inputs with differing High inputs with a “built in” oracle that fails the hypertest if the outputs are not Low equivalent. This is the method we use to detect information leaks in this paper. To repair information leaks, we need to estimate leak size.

## 2.1 Quantified Information Flow

Historically, much of the research focused on how to *check* that code obeys its designated security policy *before deployment* and was heavily influenced by Volpano, Irvine, and Smith’s work on security type systems [45]. This led to tools such as the Jif compiler and IDE [7]. One difficulty lies not in the type system approach itself, but in the noninterference property, recognized to be overly restrictive for real programs – for example, a password checking program famously will not satisfy the property for a security policy that protects correct passwords as every failure leaks a small amount of information. Attempts to ameliorate this more relevant to this paper is research in Quantified Information Flow (QIF) using information theory which had the original aim that quantitative security policies could allow information to leak, but only in a bounded way [4, 12]. While it was eventually realized that bounding QIF or exact calculation of QIF for programs is, in the worst case, computationally feasible though intractable (see Terauchi and others [48]), the attractiveness of the idea means research continues on approaches to the bounding problem [11].

We, however, aim to detect the *existence* of a leak via hypertesting, localize its cause in the code, then use Genetic Improvement (GI) to eliminate or reduce the leak size. The localization and elimination steps rely on QIF estimates but we do not need to provide guarantees for the leak bound. While it is true that eliminating leaks completely is provably equivalent to satisfying noninterference [13], as a test-based methodology, we do not guarantee the absence of leaks, only that we *may* find and fix them.

We follow our prior work [32] in using the conditional Shannon entropy of Low outputs given Low inputs as the measure of leakage, then estimating this using test sets that we assume are sampled from a discrete uniform distribution. Ultimately, the fundamental definitions were provided by Shannon [16], and the leakage

definitions by Clark et alia [13]. In what follows we present two leakage definitions and comment on how they relate to programs and security policies.

**DEFINITION 2 (LEAKAGE MEASURE FOR DETERMINISTIC PROGRAMS).** *Let  $H$  and  $H'$  be the random variables in the High parts of the initial and final states of program  $P$ , respectively while  $L$  and  $L'$  are the corresponding random variables in the Low parts. Then if  $P$  is a deterministic program and High and Low parts partition each state,  $\mathcal{L}(P)$ , the leakage from  $H$  into  $L'$  for  $P$  is given by  $\mathcal{L}(P) = \mathcal{H}(L'|L)$  where  $\mathcal{H}$  is the Shannon entropy of a random variable.*

This measures the information that flows from a random variable in the High part of the initial states to another random variable in the Low part of the final states, on the following assumptions: (1) we account for all contributions to the initial states, (2) these are partitioned between High and Low, and (3) the program is deterministic; once we account for the Low part of the initial states any remaining entropy in the Low part of the final states must be due to the High part of the initial states.

The advantage in this specialized definition is that you do not need to know anything about the random variable in the High part of the initial states. But its underlying assumptions break if some other source of information correlates with the Low part of the final state during execution. In particular, use of the definition in a test-based scenario can be sensitive to external and internal nondeterminism during executions, e.g., changing configuration parameters or race conditions in threads. In the presence of noise, Clark et al. recommend using the more general definition,  $\mathcal{L}(P) = \mathcal{I}(H; L'|L)$ , where  $\mathcal{I}$  is mutual information, and show that the two definitions are equivalent under the assumptions [13].

In our method, because we don’t need a precise bound on leakage, some flakiness in hypertests is tolerable. The tradeoffs inherent in multi-objective solutions will not always include reducing the leakage to zero. Also, using a definition that is agnostic about the entropy in the High part of states allows us to perform *experiments* where control and observation is partial as opposed to true testing where control and observation is complete. A common security policy for Unix utilities is to label inputs and outputs Low and designate data in the memory of the process as High [28]. We use the hypertest set to estimate leakage. A conventional test set could be used but a hypertest set is more focused on potential leakage.

While much software is open source, security policies tend to be implicit at best and certainly not open source. We thus re-use ones from previous work [32], outlining others in Section 4.1.

## 2.2 Related Work

In our prior work [32] we proposed the use of hyper-testing and genetic improvement to detect and fix information leaks. We demonstrated the initial promise of this approach on a small set of programs. However, by using a single-objective approach, our method necessarily folded tradeoff decisions between leakage and functionality into our chosen fitness function. We believe that this tradeoff should be available to developers and consequently that leakage repair should be based on multi-objective search. In order to achieve good repair results we reported user input was necessary. In particular, the user had to add ingredients (e.g., variables) to the search engine to help find a repair. In this work we have automated all

the steps that previously required user input. The user of LeakReducer is now only required to provide a program, its security policy, and two test suites (see Sections 3.2 and 3.1).

Some interest in hypertexting programs for information leaks has been evident in the fuzzing research community. Rather than focus on semantic leaks as we do, the interest has been in side channel leaks. CT-fuzz [22] and QFuzz [35] are examples of recent research in this area. Neither deals with automating repair. CT-Fuzz looks for failing hypertexts as evidence of leakage. Its oracle is coarse, with an observation power limited to path divergence but including timing differences. QFuzz uses a leakage measure based on min entropy to analyse the size of leaks from timing channels. Our work relates to that of the automated program repair community, but our methodology is very different. There does exist work on automated program repair in the side channel leakage community, for example, Athanasiou et alia's work [8]. This work ultimately derives from Agat's work on masking timing channels [3], but is more sophisticated, exploiting creation of statistical independence in the representation of secret data to do the masking.

### 3 LEAKREDUCER

We now present our framework for multi-objective information leakage reduction, LeakReducer. Figure 1 shows an overview.

First, a user needs to provide a program and a security policy (stage (a) in Figure 1). Next, LeakReducer requires two test sets: (1) a hypertext set, described in Section 3.2, stage (b) in Figure 1; and (2) a functional test set, described in Section 3.1, stage (c) in Figure 1. The functional test set is used to check whether current functionality is preserved; and the hypertext set is used to quantify information leakage. Assuming the hypertexts discover a non-zero information leak, then, the repair stage, stage (d) in Figure 1, can start. This stage follows the usual genetic improvement process (Section 3.3): localization of the most promising parts of program for optimization (in our case leak reduction, described in Section 3.3.2); generation of candidate patches, (Section 3.3.3); and search over the generated patches guided by a fitness, (Section 3.3.4). Finally, we discuss tool integration details in Section 3.4.

To recap, the only manual steps required are: provision of a program with its security policy, and two test sets. LeakReducer automates leak localization, detection, and repair.

#### 3.1 Automated Functional Test Generation

We automatically generate test sets that capture the current functionality of the given software. Any existing test suite could be used, but such tests are not always available. For instance, no tests covered the faulty function in OpenSSL before the Heartbleed bug was discovered. We always assume a real-world scenario where it's simply unknown whether a bug exists or not and discover any problems using hypertexts, in contradiction to, say, Mehtaev et al. [30] where their automatic repair of the HeartBleed bug was reliant on test cases only added *after* the bug was discovered.

There is a plethora of automated test generation techniques to choose from [6, 20, 46]. We decided to use fuzzers due to their effectiveness at finding software vulnerabilities [29] and the existence of multiple open source fuzzers for C which is the language

of our experimental subjects. We leave alternative functional test generation approaches as future work.

We have introduced two improvements to fuzzing described next. Intuitively, normal use of a given program is unlikely to trigger information leaks (assuming the developers did their due diligence). It is the rare, perhaps malicious, inputs that are likely to reveal such faults. Grey-box fuzzers generate inputs, driven by the goal of maximizing code coverage, and might sometimes miss such rare events. Therefore, we use a program transformation technique, HashFuzz [31], in conjunction with a fuzzer to increase the diversity of generated inputs for one of our test generation techniques. Given the rare nature of leakage triggering paths, and inability to rely on coverage feedback for discovering leaks, we felt that increased input diversity could be key to automated discovery of future leaks.

Moreover, we noticed that even with HashFuzz certain branches can be missed. As an example, Atalk [23] (one of our subject programs) has a branch that is only reachable by matching a 32-bit variable with a specific constant: `TCP_ ESTABLISHED`. If the fuzzer was only generating values for this particular 32-bit variable, then the probability of discovering this branch for each generated input would be just  $1/2^{32}$ . To further increase code coverage, we allow for the user to provide input seeds as a basis to begin fuzzing.

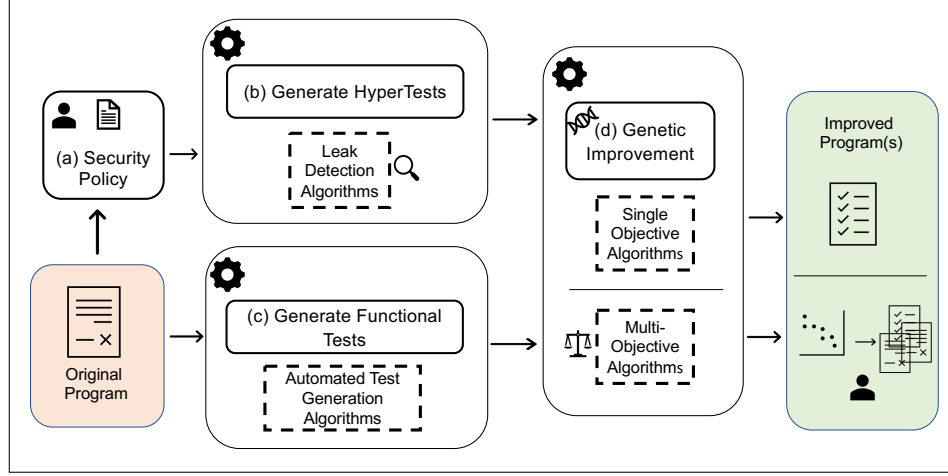
#### 3.2 Automated Hypertext Generation

As shown in HyperGI [32], traditional tests are not always able to reveal information leaks. Therefore, *hypertexts* are needed to detect and measure the amount of leakage from the target program. We use a strategy from our previous work [32] to automatically generate these. We prepare a set of Low and High inputs based on the security policy. The user provides 1) the security policy mapping, 2) the number of Low values ( $K$ ) to use in the hypertexts ( $K = 50$  in our case), and (3) the number of High values ( $N = 5$ ) to use. In the case where the High value is an input to the program, we generate  $N$  High inputs. If the High value is internal we simply run the test  $N$  times using the same Low value.

Our algorithm uses a nested binary search and a priority queue of size  $K$ , keeping the  $K$  hypertext sets with the largest quantified information flow (QIF). These sets are ultimately merged (removing duplicate pairs) into a single, large, hypertext set which will be used on the program during repair. The algorithm halves the Low input space at each iteration (if we have a 32 bit Integer type then there are 32 possible ranges), until it cannot be divided any further.

For each half Low inputs are randomly selected from within its range using  $K$  in combination with an amplifier variable  $M$  (in our case  $M = 4$ ).  $M$  is used to ensure we generate a sufficient number of values since the search spaces can be very large. We end up with approximately 200 Low inputs at the end.

For each of the Low inputs, if High inputs are needed, an inner binary search repeatedly divides the High input space in a similar fashion.  $N$  High inputs are randomly generated from each half. The hypertext sets are then run on the target program and the QIF is calculated. The hypertext sets are added to a priority queue; those with the smallest QIF are dropped. This process repeats until the binary search on the Low input space is complete. The tests are then merged and the result is a single hypertext set optimized for leakage detection.



**Figure 1: Overview of LeakReducer.** The starting point is a potentially leaky program. First a security policy is used to generate hypertests for leak detection. Automated test generation is used for functional tests. Using these test suites the program is improved using either a single or multi-objective algorithm. The result is either a single program or a Pareto front (PF) of programs from which the developer can choose the preferred patch.

### 3.3 Genetic Improvement

Once LeakReducer has generated functional (Section 3.1) and hypertest (Section 3.2) test sets, we can define the fitness functions which will guide the search algorithms in the leak reduction stage (stage (d) in Figure 1). This stage utilizes Genetic Improvement (GI) [36] to find improved program versions. The following subsections describe each of the steps of the process.

**3.3.1 Fitness Function.** The aim of LeakReducer is to find a patch that will reduce information leakage, whilst preserving software functionality. Since the two objectives can be in conflict with each other (albeit not always), we need to quantify them separately.

The first objective is quantified by the fail rate of the functional tests. That is, all functional tests should pass when run on the unmodified program, as they are used as a proxy for the intended program behaviour.

The second objective is quantified by using a hypertest set to estimate the leak size. We use the same calculation as in previous work [32] (see Definition 2 in Section 2), and henceforth use QIF to refer to estimate for quantified information leakage.

**3.3.2 Leak Localization.** Before starting the search process, possible leak locations need to be identified. Similar to previous work [32], we use a lightweight form of dynamic analysis for this purpose, albeit on the whole leaky file not just the function. Based on the security policy, we first identify the file where leakage occurs. Next, we remove each statement from the target, one-by-one, and observe the impact on QIF when the hypertests are run.

When a statement is removed, if the program fails to compile or run, then QIF for that statement is assumed to be zero. Otherwise, we calculate and store the absolute change in QIF of this modified program from the original. Finally, statement leakages are normalized with respect to the highest absolute difference. The higher the absolute difference, the greater the probability that the fault is

coupled with that statement. We use these probabilities during the search process to prioritize statements that have higher impact on information leakage based on this analysis. The assumption is that changes to those statements should have highest impact on (and hopefully reduce) the leakage.

**3.3.3 Patch Generation.** Once we have identified statements that have influence on information leakage, we mutate them to create new program variants. In genetic improvement the standard mutation operators simply insert, delete and replace software fragments, e.g., code statements. We observe that information leakage problems are caused by information flow between Low and High variables in a given program. Therefore, aside from traditional GI operators, we use the *NewIf* and *NewFor* mutation operators, introduced in our previous work [32]. In LeakReducer, we fully automate this step of the process. In order to create a new statement, we need expressions to populate the parentheses for FOR and IF statements. To reduce the search space, we construct comparison expressions using existing identifiers and type match them, so that, e.g., a number is compared with another number. For both *NewIf* and *NewFor*, the body statement to be executed is either selected from existing code or one of four code fragments, 1) *return 0*; 2) *return 1*; 3) *id1 = id2*; or 4) *id2 = id1*; are inserted.

**3.3.4 Search Strategies.** Another improvement we introduce over our previous work [32] is a choice between search strategies, including multi-objective search. We aim to reduce both 1) the number of failing test cases, i.e., *fail\_rate* and 2) information flow leakage, measured with *QIF*. The optimal solution will reduce both of these quantities to 0.

**3.3.5 Single Objective Optimization.** The single-objective search option uses the same fitness function as previously [32].

$$fitness = (fail\_rate + \frac{QIF_i}{QIF_0})/2 \quad (1)$$

As can be seen from Equation 1, this fitness function balances the two objectives: *fail\_rate* and normalized *QIF* where  $QIF_0$  is the initial leakage ( $QIF_0 > 0$ ), and  $QIF_i$  is the leakage of current program variant ( $QIF_i \geq 0$ ).

**3.3.6 Multi-objective Optimization.** In multi-objective optimization, there are several objectives where maximizing (or minimizing) one objective may have conflicting results with other objectives. LeakReducer has two objectives: reducing the information leakage and preserving intended functionality, which can be in conflict. An example is the Triangle program, presented in our previous work [32], where reduction of information leakage is impossible without loss of triangle classification correctness.

LeakReducer’s multi-objective [15] setting thus reports a list of non-dominated Pareto front solutions, where every solution in the list is better than the other solutions in at least one objective. This provides the decision maker with more options to choose from.

## 3.4 Integration

In order to implement LeakReducer, we extended an existing genetic improvement framework, PyGGI [5] and integrated it with the JMetalPy framework for multi-objective optimization [10]. This integration was a non-trivial task, as it required integration of frameworks with different architectures. PyGGI provides us with the GI framework, while JMetalPy provides different single and multi-objective operators, algorithms and Pareto front quality metrics.

For functional test case generation, we investigated several state-of-the-art fuzzing strategies (Section 3.1). We use AFL++ 3.12 [19] which consistently ranks amongst the best fuzzers in terms of program exploration ability [33]. The HashFuzz transformation was also applied to the tested programs, and these were fuzzed using the same setup as the untransformed programs. We also tested manual seeds for both fuzzing variants, to increase path coverage (we refer to this as Test Augmentation - abbreviated to TA). Thus we had four set-ups: AFL, HashFuzz, AFL+TA, and HashFuzz+TA.

Lastly, we built an automated identifier extractor for patch generation (Section 3.3.3). We use *Universal Ctags* [44] (or *ctags* for short) to extract an initial set of variables with their types from the target files. To enrich ingredients, we exploit PyGGI’s internal program representation (tagged XML) to extract expressions. We then use a custom function that infers their type. For example, assume we have the following statement: `int len = 2 * y, dist;` *ctags* detects `len` and `dist` and deduces that these are both integers. From this, we can infer that the expression `2 * y` can be used as an integer.

## 4 EVALUATION

Our aim is to provide an effective automated tool that can estimate and reduce information leaks in real-world software. We therefore ask three research questions to evaluate LeakReducer. Our first question focuses on detection of information leakage:

**RQ1: How effective is LeakReducer at detecting information flow leakage in a given program?**

We compare LeakReducer’s leak detection algorithm against a state-of-the-art fuzzing tool, using different test seeding strategies.

Our next RQ focuses on comparison with a single-objective search strategy, originally proposed in our previous work [32]:

**RQ2: How effective is LeakReducer at reducing information flow leakage using a single-objective algorithm?**

Even if LeakReducer is successful at reducing leakage using single-objective search, we still need to answer our key hypothesis: that we need to balance leakage and functionality. Thus we ask:

**RQ3: How effective is LeakReducer at reducing information flow leakage using a multi-objective algorithm?**

We use multiple multi-objective algorithms and examine the results from a qualitative perspective, to see if several viable solutions are found on the Pareto front.

### 4.1 Benchmarks

We use six programs for our study, with associated security policies. We present these in Table 1. The first three programs are the same that were used in previous work [32]. Two of these have been used in prior research on information leakage [23] and reported in the Common Vulnerabilities and Exposures (CVE) database [40].

The Classify program is similar to the Triangle program, in that it outputs a class depending on the input parameters. In contrast to the Triangle program, however, the output space is further divided – rather than 3 possible outputs (scalene, isosceles and equilateral) in Triangle, there are 11 in Classify. We introduced this program to specifically show tradeoffs between leakage and functionality loss.

The last two programs are real programs meant to show the scalability and practicability of LeakReducer. Both are taken from the OpenSSL library [41, 42]. The first is: `dtls1_process_heartbeat` function from `openssl-1.0.1f` containing the Heartbleed bug, which leaks information from internal memory.

For our second program, we chose another OpenSSL library, but one that has no known leaks. We examined functions within the OpenSSL crypto directory, to find one which might have interactions between Low and High security information and with a similar signature to Heartbleed. In the same version of OpenSSL we found the `BIGNUM *BN_bin2bn()` function in the `bn_lib.c` file of the BigNum library. This function is one of the data entry functions for the library. Given the similarity between Heartbleed and BigNum, security policy creation was straightforward (see Table 1). Although, we used `openssl-1.1.1j` [43] in our tests, the same function is still in use in up-to-date OpenSSL versions.

### 4.2 Experimental Protocol

Next, we describe the methodology we use to answer our RQs.

**4.2.1 Automated Test Generation.** To answer RQ1 about LeakReducer’s effectiveness in finding information leaks, we automatically generate two types of tests: functional (Section 3.1) and hyper-tests (Section 3.2). Although we use functional tests primarily to establish intended program behaviour, such tests might also reveal information leakage. In particular, we use AFL fuzzer’s address sanitizer [27] option, which allows detection of subtle memory access issues (such as out-of-bounds or use-after-free).

We use four different Fuzzer-Test Setups (FTS): 1) AFL 2) AFL-TA 3) HashFuzz and 4) HashFuzz-TA, as discussed in Section 3.4. For the Test Augmentation (TA) phases we use up to 2 manual seeds, to reach the hard to find paths that fuzzing may miss. We first run each fuzzer-test setup 5 times with the `AFL_EXIT_WHEN_DONE` option. With this option set, AFL automatically terminates when

**Table 1: Study subjects.** For each we provide: a reference; the lines of code in the file containing the function of interest (which LeakReducer targets); a CVE number if information leakage was reported for this function; and the security policy used, with parameters from the function’s signature and *function return values*.

Subject	Ref	LoC	CVE-#	Security Policy		
				High input	Low input	Low Output
Triangle	[32]	14	–	secret	side2 & side3	<i>function return value</i>
Atalk	[23]	33	CVE-2009-3002	<i>internal memory</i>	sock & peer	<i>function return value &amp; uaddr</i>
Underflow	[23]	18	CVE-2007-2875	h	ppos	<i>function return value</i>
Classify	authors	18	–	High	Low	<i>function return value</i>
Heartbleed	[41]	1,082	CVE-2014-0160	<i>internal memory</i>	payload_sent & payload_length	payload_received
Bignum	[42]	778	–	<i>internal memory</i>	s, len & ret	s & <i>function return value</i>

all discovered paths have been fuzzed many times without any new path being found [2]. From these first 5-runs, we identified the *max fuzz-time*. Then, to maximize path coverage, we run each FTS 5 times again for *max fuzz-time* seconds.

To generate hypertexts we use the algorithm described in Section 3.2. For each of the test sets we report the number of unique test cases, generation time budgets, crashes found (whether they reveal leaks or not), as well as QIF values for hypertexts, in Table 2.

**4.2.2 Parameter Tuning.** Since genetic programming has been the dominant search strategy in GI [36], we use it in LeakReducer’s single-objective setting. However, it’s not obvious which parameter settings would be optimal for our application domain. Therefore, we conduct parameter tuning on the four smaller subjects, (Triangle, Classify, Atalk and Underflow) before running the information leakage reduction stage of LeakReducer. We vary: population size – between 20 and 100, in increments of 20, as these are the min and max found in the GI literature – and mutation and crossover rates – each tried with four values: 0.25, 0.5, 0.75 and 1. We use a budget of 2,000 program evaluations for each setting. Overall, we test 80 configurations of single-objective LeakReducer for each of the 4 settings and 4 test subjects, repeated 5 times each, for a total of 6,400 individual runs.

Recall that we use functional test fail rate and QIF from hypertexting to calculate the fitness value of each program variant. Thus, in order to fairly compare runs which use different functional test sets, we combine all 4 test sets for evaluation of the best individual found in each run. This is not an issue for hypertexts, which are generated once for all subjects.

**4.2.3 Leak Reduction.** To answer RQ2 and RQ3 about LeakReducer’s effectiveness at leak reduction, we run LeakReducer’s GI stage using two search strategies: single- and multi-objective. We set the GI budget to 10,000 evaluations. We repeat each repair run 20 times. As before, the fail rate in the reported final fitness calculation is based on runs of the combined set of all functional tests.

**Single-Objective Optimization** We use the best parameter settings from the tuning stage for each of the functional test sets for each of our subject programs. We record the fitness values, runtimes, as well as functional test fail rates and QIF values. We also re-run evolved individuals from the HyperGI work.

**Multi-Objective Optimization** We selected four multi-objective optimization algorithms: NSGAII, NSGAIII, MoCell and SPEA2. We

selected the first three, as they showed good performance in recent work tackling a multi-objective problems [39]. NSGAIII [17] was originally proposed for many-objective optimization (3 or more objectives) to help scalability. However, it has been used for problems with only 2 objectives [39], hence we chose to include it in our experiments as well. NSGAIII requires reference points to be set. We used the UniformReferenceDirectionFactory from the JMetalPy framework [10] and set the number of the reference points equal to the number of the individuals in the population minus 1 which is similar to [39]. For the other parameters, we used the same settings as we used in NSGAII. We also added SPEA2, since it proved successful in improvement of non-functional properties of software using genetic improvement [47]. Based on Li et al. [26]’s guidance, we report the following quality indicators of the Pareto fronts generated by the selected algorithms: Hyper-volume, Inverted Generational Distance, Epsilon and Generational Distance.

### 4.3 Test Environment

Due to compatibility issues, we ran the Heartbleed fuzzer experiments on a single core of an Intel (R) Core(TM) i5-2500 CPU @ 3.3GHz processor, with 8GB RAM and 500GB HDD with Ubuntu 20.04 and gcc & g++ 7.3.0. All other experiments were run on a High Performance Computing (HPC) cluster using single cores of an Intel(R) Xeon(R) Gold 6244 CPU @ 3.60GHz and 8GB of RAM, in RedHat Enterprise Linux 7 with gcc & g++ 10.2.0.

### 4.4 Threats to Validity

With respect to generality, we experimented on a limited number of programs, some with similar characteristics. However, (1) we wanted to compare against the state-of-the-art and (2) security policies, which are necessary for detecting information flow leakage cannot be developed without extensive knowledge of a system. Hence, we focused on programs that have policies based on the High/Low lattice. With respect to internal validity, we acknowledge there could be faults in our programs, but we have manually validated our patches, and are providing them on our online website along with our other artifacts. We also contacted developers of the program for which LeakReducer found a possibly new information leakage fault. Lastly, with respect to construct validity, we could have chosen different metrics, but we have used the most common metrics for evaluating multi-objective optimization and use the same measure of leakage as prior work.



## 5 RESULTS

In this section, we answer each of our research questions. Artifacts for the experiments are found in our github repository [1].

### 5.1 Parameter Tuning Results

We used Wilcoxon rank sum test over fitness values from the parameter tuning tests. 15 of 80 configurations showed evidence of a difference in means ( $p - value \leq 0.05$ ) and 5 configurations had weak evidence ( $0.05 < p - value \leq 0.1$ ). Out of the 15 configurations which were significant, we chose the configuration with the smallest mean fitness and the largest population: population size = 100, mutation rate = 1, and crossover rate = 0.5. We thus use these parameters in the following experiments.<sup>2</sup>

### 5.2 RQ1: Leak Detection

To answer this question, we examine the data in Table 2. All of the fuzzing variants found crashes for Heartbleed and Bignum. We examined the crashes, and indeed they were caused by memory issues that lead to leakage. As the Bignum crash was not from a confirmed bug, we reached out to the developers who confirmed our suspicion that this particular function is typically not accessible from outside of a program, and in the case that it is, it would be the responsibility of application developers to ensure proper use of the function. However, the code has similarities to that found in many other confirmed cases of information leakage in the wild. For the other four programs no crashes were produced by fuzzing, even though the time budget for each program was significantly higher than the time for generating hypertexts. Furthermore, running LeakReducer’s hypertexts revealed leakage in all 6 subjects, producing non-zero QIF values in all runs.

**Answer to RQ1 (Leak Detection):** LeakReducer is able to detect and estimate leakage in all 6 subjects, where state-of-the-art fuzzers only find leakage-related crashes in 2 of the subjects.

### 5.3 RQ2: Single-objective Leak Reduction

To answer RQ2 we turn to Table 3. It shows the normalized post patch (i.e., improved) QIF (QIF<sub>i</sub>), Fail Rate (FR) and Fitness (F) for HyperGI and single-objective LeakReducer with 4 different inputs from fuzzing. To compare LeakReducer’s results with results obtained by Mesecan et al. [32], we ran each evolved program using our test cases (all 4 functional test sets + hypertexts). The best fitness values are highlighted in bold. With respect to previous work [32]: 1) we obtained the same median fitness values with the AFL-TA setting for Underflow; 2) we obtained worse median fitness values for Atalk; 3) and, we obtained better median fitness values for the Triangle program for all fuzzer test setups.

Interestingly, LeakReducer did not find the optimal solution for Atalk. Further investigation revealed this is because in [32] we allowed the user to input extra identifiers for mutation creation in order to alleviate a weakness of our previous implementation, which could only extract and detect types for a limited set of identifiers. The results in our previous work depended on developer expertise. In LeakReducer, we improved the identifier detection

<sup>2</sup>The actual mutation rate is closer to 0.5 based on JMetalPy’s implementation.

module by expanding automated extraction and inference of identifier types. This is now fully automated, though it came at a cost of a larger search space, and restriction of identifiers to the target file only. The identifier provided by a developer in [32] was not used in the target file, i.e., LeakReducer had no access to it and a solution that fully reduces leakage was not found. We argue, however, the automation of the process outweighs potentially missing such fixes. LeakReducer still found some program variants that reduced leakage. These partial fixes could still be helpful for the developers for understanding and fixing such leaks. Moreover, in the future, the ingredient space can be enlarged, e.g., by using new mutation operators, such as a memory initialization mutation operator.

To choose the best test-generation setup, we used the Wilcoxon rank sum test on the fitness values. AFL-TA provided slightly better median fitness for 5 out of 6 test subjects. But, it had a  $p - value = 0.102$ , hence there is no statistical difference between AFL-TA and the other fuzzer test setups. Similarly, the Wilcoxon rank sum test between Automated Test Generation (ATG) and Test Augmentation (TA) showed no significant difference ( $p - value = 0.101$ ).

**Answer to RQ2 (Single-objective Leak Reduction):** Single-objective LeakReducer achieves results competitive with previous work [32], but without the need for human interaction during the repair process. Search for repairs guided by tests from AFL with seeded inputs (AFL-TA) did produce marginally better results than other settings, but the result was not statistically significant.

### 5.4 RQ3: Multi-objective Leak Reduction

To compare different multi-objective variants, we used 4 quality indicators, as previously discussed. Since we did not see a clear winner for the fuzzer test setups and to ensure high coverage, we simply used all functional tests generated from the previous steps, and used those during search to evaluate correctness of the evolved program variants. To compare Pareto fronts from different runs, we first prepared the *global Pareto front*, i.e., Pareto front that combines all generated fronts. Next, we calculated the hypervolume and distance from the global Pareto front to individual fronts. Table 4 reports median distances and hypervolume for each subject and algorithm. Counting the number of times each algorithm produces the best result (largest Hypervolume or Minimum Median Distance to the global Pareto front for the other metrics), we get: (1) SPEA2: 19; (2) NSGA-III: 15; (3) NSGA-II: 11; and (4) MOCcell: 8. These results indicate that the SPEA2 option performs competitively well in this problem domain, with NSGA-III not far behind. Consequently, we now look more closely at the Pareto fronts generated by LeakReducer using SPEA2.

Figure 2 shows results of individual runs for SPEA2 for 3 test subjects and Figure 3 shows the Pareto fronts from all runs for all subjects for SPEA2 together with the best solution from single-objective LeakReducer marked.

Interestingly, the single-objective variant produces results that are on or close to multi-objective LeakReducer’s Pareto fronts. However, the latter option provides more variants to choose from, leveraging leakage vs fail rate. Moreover, whether improvement of information leakage conflicts with intended program semantics will



**Table 2: Leak detection and functional test case generation using four variants of fuzzing compared against LeakReducer’s Hypertest set. For each of the fuzzing settings we show if the system crashed. For LeakReducer’s Hypertest set we show the original QIF<sub>0</sub>. We also give the test suite size (TS), max fuzzing time budget and hypertest generation runtime in hours.**

Subject	AFL			AFL-TA			HashFuzz			HashFuzz-TA			Hypertests		
	Crashes	TS	Hours	Crashes	TS	Hours	Crashes	TS	Hours	Crashes	TS	Hours	QIF <sub>0</sub>	TS	Hours
Atalk	No	12	1.0	No	18	0.9	No	122	16.2	No	116	12.9	1.1	50	1.0
Bignum	Yes	3	7.8	Yes	24	6.4	Yes	12	21.2	Yes	70	34.7	2.4	40	0.1
Classify	No	107	2.6	No	122	5.7	No	178	8.5	No	176	8.1	2.4	125	0.2
Heartbleed	Yes	47	27.8	Yes	60	76.1	Yes	84	266.1	Yes	109	145.4	2.3	40	0.4
Triangle	No	60	3.1	No	54	1.7	No	66	58.3	No	78	69.2	0.8	194	0.1
Underflow	No	20	0.4	No	40	0.5	No	234	13.4	No	252	16.2	5.4	100	1.3

**Table 3: Fitness values for best evolved variants in Mesecan et al.’s work [32] and single-objective LeakReducer, guided by 4 fuzzer test sets. Medians of 20 runs are shown. In addition, we report the two fitness components: fail rates (FR) and leakage estimates (QIF<sub>i</sub>). The best fitness values are highlighted in bold.**

Subject	Mesecan et al.[32]			AFL			AFL-TA			HashFuzz			HashFuzz-TA		
	QIF <sub>i</sub>	FR	Fitness	QIF <sub>i</sub>	FR	Fitness	QIF <sub>i</sub>	FR	Fitness	QIF <sub>i</sub>	FR	Fitness	QIF <sub>i</sub>	FR	Fitness
Atalk	0.00	0.00	<b>0.00</b>	1.00	0.00	<b>0.50</b>	0.51	0.24	<b>0.50</b>	1.00	0.00	<b>0.50</b>	1.00	0.00	<b>0.50</b>
Triangle	0.00	0.69	0.34	0.23	0.23	0.24	0.15	0.38	0.23	0.10	0.30	<b>0.22</b>	0.11	0.29	0.23
Underflow	0.00	0.62	<b>0.31</b>	0.11	0.57	0.35	0.11	0.56	<b>0.31</b>	0.64	0.01	0.37	1.00	0.00	0.50
Bignum	-	-	-	0.00	0.78	0.39	0.00	0.60	<b>0.30</b>	0.00	0.78	0.39	0.33	0.47	0.39
Classify	-	-	-	0.00	0.64	0.32	0.02	0.61	<b>0.31</b>	0.00	0.64	0.32	0.00	0.63	<b>0.31</b>
Heartbleed	-	-	-	0.50	0.00	<b>0.25</b>	0.50	0.00	<b>0.25</b>	0.50	0.00	0.43	0.50	0.00	<b>0.25</b>

**Table 4: Pareto front quality indicators per subject & multi-objective algorithm. Median of 20 runs reported. A larger hypervolume indicates a better result. For other indicators smaller numbers are closer to the global front, hence better.**

	Inv. Generational Distance				Generational Distance				Hyper Volume				Epsilon			
	MOCeII	NSGAII	NSGAIII	SPEA2	MOCeII	NSGAII	NSGAIII	SPEA2	MOCeII	NSGAII	NSGAIII	SPEA2	MOCeII	NSGAII	NSGAIII	SPEA2
Atalk	0.15	0.43	0.15	<b>0.14</b>	0.04	<b>0.03</b>	0.12	0.11	0.67	0.23	0.71	<b>0.72</b>	0.24	0.73	0.22	<b>0.20</b>
Bignum	<b>0.38</b>	<b>0.38</b>	<b>0.38</b>	<b>0.38</b>	<b>0.33</b>	0.37	0.38	0.35	0.46	0.43	0.43	<b>0.49</b>	<b>0.40</b>	<b>0.40</b>	<b>0.40</b>	<b>0.40</b>
Classify	0.08	0.08	<b>0.07</b>	<b>0.07</b>	0.08	0.08	<b>0.07</b>	<b>0.07</b>	0.68	0.67	<b>0.69</b>	<b>0.69</b>	0.13	<b>0.12</b>	0.13	<b>0.12</b>
Heartbleed	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>	<b>0.02</b>	0.50	0.50	0.50	<b>0.49</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>0.97</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>
Triangle	0.04	<b>0.03</b>	<b>0.03</b>	<b>0.03</b>	0.03	0.02	0.02	<b>0.01</b>	0.75	<b>0.76</b>	<b>0.76</b>	<b>0.76</b>	0.09	<b>0.05</b>	<b>0.05</b>	<b>0.05</b>
Underflow	<b>0.20</b>	0.30	<b>0.20</b>	0.58	0.24	<b>0.00</b>	<b>0.00</b>	<b>0.00</b>	0.66	0.54	<b>0.67</b>	0.00	<b>0.28</b>	0.43	<b>0.28</b>	1.00

**Table 5: Means and standard deviations of the numbers of solutions in Pareto fronts of multi-objective LeakReducer. The best values are highlighted in bold.**

	MOCeII		NSGAII		NSGAIII		SPEA2	
	Mean	Stdev	Mean	Stdev	Mean	Stdev	Mean	Stdev
Atalk	5.4	1.4	5.9	1.4	5.5	1.5	<b>6.7</b>	2.2
Bignum	<b>3.0</b>	0.0	<b>3.0</b>	0.3	<b>3.0</b>	0.2	<b>3.0</b>	0.0
Classify	<b>35.8</b>	4.3	33.8	3.7	31.9	4.1	31.7	5.1
Heartbleed	<b>2.1</b>	0.2	<b>2.1</b>	0.2	2.0	0.0	2.0	0.0
Triangle	6.9	1.8	7.3	1.5	<b>7.5</b>	1.5	7.4	1.6
Underflow	3.5	1.1	3.5	1.2	3.5	1.2	<b>4.1</b>	1.3
Average	9.42	-	9.24	-	8.88	-	9.13	-

depend on the program at hand. As mentioned before, it is impossible to reduce leaks in the Triangle program without sacrificing functional correctness, while the optimal Heartbleed information leakage fix does not change the intended behavior of the program. Therefore, in the cases where reduction of information leakage is a priority we recommend the use of the multi-objective approach.

We are also interested in the diversity of solutions found. Table 5 presents the average number of solutions found and the standard deviations out of 20 runs for each algorithm and subject pair. The largest number of solutions for each subject is highlighted in bold. Firstly, the average number of solutions in Pareto fronts is quite close; ranging from 8.9 for NSGAIII to 9.4 for MOCeII. Secondly, MOCeII and SPEA2 generated the largest number of solutions for 3 subjects, while NSGAII and NSGAIII generated the largest number of solutions for 2 subjects. We will discuss the quality of solutions in

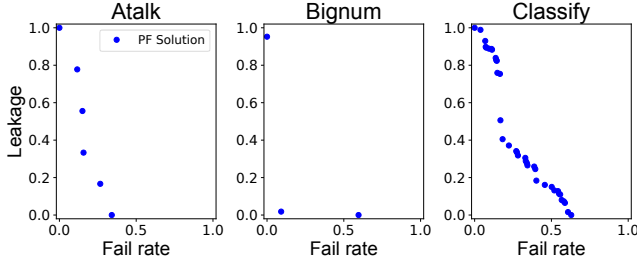


Figure 2: Sample Pareto fronts from LeakReducer with the SPEA2 setting. Each dot represents a solution on the front.

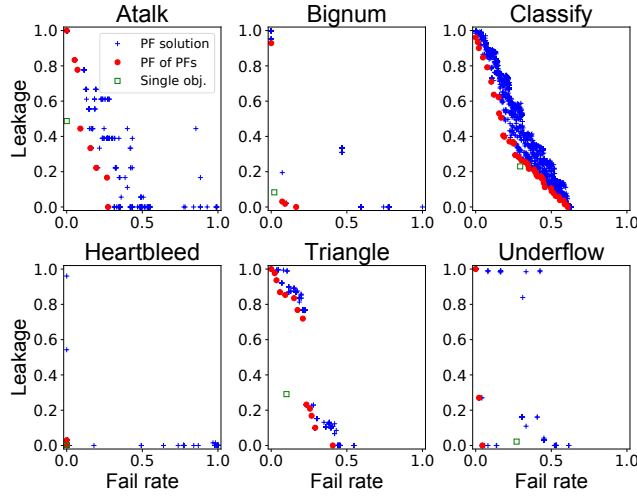


Figure 3: Pareto fronts from all 20 LeakReducer runs with the SPEA2 setting (blue), and Pareto front of all those 20 fronts (red). The green square shows the best solutions found by single-objective LeakReducer.

the next section. Taking all results into account, however, we would recommend the use of SPEA2 with LeakReducer, as it is likely to produce a diverse set of solutions, close to the optimal Pareto front.

**Answer to RQ3 (Multi-Objective LeakReducer):** Both single and multi-objective settings produced effective fixes, with single-objective LeakReducer producing solutions close to the multi-objective LeakReducer’s Pareto fronts, while the latter provided a diverse set of solutions balancing leakage and fail rate. Moreover, we recommend the use of SPEA2 with LeakReducer, as we showed with multiple indicators that it generally outperforms other multi-objective algorithms.

## 5.5 Discussion

We now discuss some of the patches obtained by LeakReducer and our observations in more detail.

**5.5.1 Atalk Leakage.** Atalk is the real-world program that showed the most evenly spread Pareto front in Figure 2. The original leaky function is:

```

struct atalk_sock {
    unsigned char dst_node, src_node, dst_port, src_port;
    int sk_state;
    char res[16];
};

int atalk_getname(atalk_sock *sock, atalk_sock *uaddr, int peer) {
    struct atalk_sock sat;
    int err = -ENOBUFF;
    if (sock_flag(sock))
        goto out;

    if (peer) {
        err = -ENOTCON;
        if (sock->sk_state != TCP_ESTABLISHED)
            goto out;
        sat.src_node = sock->dst_node;
        sat.src_port = sock->dst_port;
        sat.dst_node = sock->src_node;
        sat.dst_port = sock->src_port;
    } else {
        sat.src_node = sock->src_node;
        sat.src_port = sock->src_port;
        sat.dst_node = sock->dst_node;
        sat.dst_port = sock->dst_port;
    }

    sat.sk_state = sock->sk_state;
    memcpy(uaddr, &sat, sizeof(sat));
    err = sizeof(atalk_sock);

out:
    return err;
}

```

This particular program leaks values from internal memory due to the struct `sat` defined on line 2 being uninitialized. There are 6 struct members, but only 5 are assigned to in the code in lines 3-22. The entire memory contents of the struct are then copied to the function parameter `uaddr`, including the value of the uninitialized 6th struct member `res`. As `sat` is a local variable, this uninitialized memory contains stack data, which could, depending on previous function call stacks, contain sensitive data.

**5.5.2 Atalk Patches.** For the sake of brevity, we show 2 of the best patches generated by SPEA2:

Patch 2; Leakage: 0.000, Functional: 0.341

```

10a11
> sat.src_port = sock->src_port;
18a20,22
> if (sock->src_port < sock->src_node) {
>     sock->src_node = sock->src_port;
> }
20a25,27
> if (sock->dst_port != sock->src_node) {
>     return 1;
> }
23a31,33
> for (int lcv1476 = 0; lcv1476 < peer; lcv1476++) {
>     err = -ENOTCON;
> }

```

Patch 5; Leakage: 0.778, Functional: 0.116

```

6a7
> err = -ENOTCON;
10a12,14
> for (int lcv943 = 0; lcv943 < TCP_ESTABLISHED; lcv943++) {
>     sat.src_port = sock->src_port;
> }
21a26,28
> if (sizeof(sat) > sock->src_port) {
>     return 0;
> }

```

There is variation in these patches, and they are both returning a fixed value before the leak occurs (in line 23) depending on a comparison between structs `sat` and `sock`. The variations in leakage and functional test performance can be attributed to different comparisons between the two structs. These comparisons cause an early return before the leakage occurs in differing proportions of the test inputs, hence differing functional results and leakage rates.

An ideal patch would initialize the `sat.res` struct member to a fixed value. As this variable is in fact an array, this would require a `for` loop, or a call to `memset`. As an alternative, using the short form compound literal initializer `struct atalk_sock sat = {0}`; will initialize all struct members (including all array members) to all zeroes. Neither of these are used in the existing code, so the GI would require additional mutation operators to produce these.

The 3 other real-world programs (Heartbleed, Bignum and Underflow) all contain developer bugs, which leads to a less diverse range of solutions. They can in theory all be ‘fixed’ in a way that retains all original functionality, whilst completely eliminating the leakage. Both Classify and Triangle have leakage that is caused by intentionally poor information flow control design, as such these produce a broad, dense Pareto front showcasing a variety of potential solutions.

**5.5.3 Heartbleed Leakage and Patches.** The best single-objective patch found for Heartbleed was as follows (purple code was removed in the patch, and green code was added by the patch):

```

1  int dtls1_process_heartbeat(SSL *s) {
2      ...
3      unsigned int payload, padding = 16;
4      n2s(p, payload); /* Read payload length */
5      ...
6      unsigned char *buffer = OPENSSL_malloc(1 + 2 + payload + padding);
7      unsigned char *bp = buffer;
8
9      *bp++ = TLS1_HB_RESPONSE; /* Copy response type into buffer */
10     s2n(payload, bp); /* Copy payload length into buffer */
11     memcpy(bp, p, payload); /* Copy payload into buffer */
12     bp += payload; /* Update buffer pointer */
13     + if (SSL_F_DTLS1_PREPROCESS_FRAGMENT < 1 + 2 + payload + padding) {
14     +     return 0;
15     + }
16     - RAND_pseudo_bytes(bp, padding);
17
18     /* Send callback */
19     r = dtls1_write_bytes(s, TLS1_RT_HEARTBEAT, buffer, 3 + payload + padding);
20     ...

```

The added `if`-statement (lines 13-15) fixes many leaks because of the way that Heartbleed is exploitable. Essentially the payload length in the malformed packet is set to some value larger than the length of the actual message sent to be echoed back. padding is always 16. Multi-objective search results include patches that are semantically equivalent to the single-objective patch shown above. The payload length is read in by a function-like macro `s2n` on line 10, which reads an unsigned 2-byte (16-bit) integer (between 0-65535) from the transmitted heartbeat buffer. The leak occurs when payload bytes are copied into the callback message on line 11 and sent back to the other client; thus in order to maximize leakage, an attacker should set this value to 65535 (the maximum) which will result in  $65535 - 19 - \text{ACTUAL\_PAYLOAD\_LENGTH}$  bytes of internal program memory being returned to them in the corresponding heartbeat response. The `-19` comes from the 16 bytes of padding and an extra 3 bytes for indicating message type and payload length. The constant `SSL_F_DTLS1_PREPROCESS_FRAGMENT` is defined as 288 in a header file. This patch is therefore discarding any heartbeat

request where payload is larger than 269 ( $288 - 19$ ), which reduces the leakage from a potential 65516 bytes down to 269 bytes. An attacker looking to maximize leakage would be setting payload to very large values, and this patch would silently discard these. The fuzzer however is not using any malicious heuristics when generating malformed packets, but instead generates the buffer (containing the payload length) pseudo-randomly, and with the leaking search space shrinking down to 0.41% ( $269/65535$ ) of what it originally was, the leak is seemingly repaired. It is worth noting that the patch will also discard any properly formed heartbeat packets with an actual payload of length 270 or greater.

The developer patch is close to this. The only difference is that the comparison value `SSL_F_DTLS1_PREPROCESS_FRAGMENT` is replaced with the actual length of the received buffer. This eliminates the leakage for requests with indicated payloads  $\leq 269$  bytes, whilst accepting properly formed packets with an actual payload length  $> 269$  bytes. An improved test set would expose the remaining (much smaller) leak, and is a target for future research.

The evaluation criteria does well to guide the repair process, as is evidenced by the proposed patch being very close to the actual developer patch. An acclimatized OpenSSL developer would see that comparing the variable payload length with a constant value would create issues for properly formed long heartbeat requests, but learn that the leakage is strongly correlated to large payload values. It is not a leap to suggest that a developer could come to the correct conclusion that the solution is to discard requests where the actual buffer length does not match the indicated buffer length.

**5.5.4 Bignum Patches.** Bignum leaks information via a buffer-overflow, and does not have a single ‘fix’, however, we do see patches that manage to greatly reduce the quantity of leakage while not resulting in a large functionality testcase fail rate increase. The best patches generated through multi-objective LeakReducer for Underflow reduce leakage to a greater extent than the best single objective patch, whilst simultaneously retaining greater functionality.

## 6 CONCLUSIONS AND FUTURE WORK

We have presented LeakReducer, a multi-objective framework to detect, localize and repair information leakage in real-world programs. We have evaluated LeakReducer on a set of six programs and were able to find leaks in all of them. We demonstrated our ability to repair the leaks against repairs from the state-of-the-art tool, introduced in our previous work. For those programs whose leakage and functionality are competing, the multi-objective setting provided a diverse Pareto front which can be used to balance leakage and functionality. We make artifacts available on our github repository [1] to facilitate further research in this exciting area. We plan to add additional repair ingredients, refined identifier resolution, specialized mutation operators to zero out memory and perform array bound checks, in addition to applying LeakReducer to a larger set of programs and varying security policies.

## ACKNOWLEDGMENTS

This work was supported in part by NSF grant CCF-1909688 and by UKRI EPSRC grants EP/P023991/1 and EP/P005888/1.

## REFERENCES

- [1] 2022. <https://github.com/LavaOps/LeakReducer/>.
- [2] AFL 2022. American Fuzzy Lop plus plus (AFL++). <https://github.com/AFLplusplus/AFLplusplus>. Accessed: 2022-05-22.
- [3] Johan Agat. 2020. Transforming Out Timing Leaks. In *POPL 2000, Proceedings of the 27th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Boston, Massachusetts, USA, January 19-21, 2000*, Mark N. Wegman and Thomas W. Reps (Eds.). ACM.
- [4] M. S. Alvim, K. Chatzikokolakis, A. McIver, C. Morgan, C. Palamidessi, and G. Smith. 2020. *The Science of Quantitative Information Flow*. Springer.
- [5] Gabin An, Aymeric Blot, Justyna Petke, and Shin Yoo. 2019. PyGGI 2.0: Language Independent Genetic Improvement Framework. In *Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 1100–1104.
- [6] Andrea Arcuri. 2019. RESTful API automated test case generation with EvoMaster. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 28, 1 (2019), 1–37.
- [7] Owen Arden, Jed Liu, Tom Magrino, and Andrew Myers. 2016. Jif: Java with Information Flow. <https://www.cs.cornell.edu/jif/>. Accessed: 2022-05-02.
- [8] Konstantinos Athanasiou, Thomas Wahl, A. Adam Ding, and Yunsu Fei. 2020. Automatic Detection and Repair of Transition-Based Leakage in Software Binaries. In *Software Verification - 12th International Conference, VSTTE 2020, and 13th International Workshop, NSV 2020, July 20-21, 2020, Revised Selected Papers (Lecture Notes in Computer Science, Vol. 12549)*, Maria Christakis, Nadia Polikarpova, Parasara Sridhar Duggirala, and Peter Schrammel (Eds.). Springer, 50–67.
- [9] D. Elliott Bell and Leonard J. LaPadula. 1973. Secure Computer Systems: Mathematical Foundations. MITRE Technical Report 2547, Volume 1.
- [10] Antonio Benitez-Hidalgo, Antonio J. Nebro, Jose Garcia-Nieto, Izaskun Oregi, and Javier Del Ser. 2019. jMetalPy: A Python framework for multi-objective optimization with metaheuristics. *Swarm and Evolutionary Computation* 51 (2019), 100598.
- [11] Fabrizio Biondi, Michael A. Enescu, Annelie Heuser, Axel Legay, Kuldeep S. Meel, and Jean Quilbeuf. 2018. Scalable Approximation of Quantitative Information Flow in Programs. In *Verification, Model Checking, and Abstract Interpretation - 19th International Conference, VMCAI, Isil Dillig and Jens Palsberg (Eds.)*. Springer.
- [12] David Clark, Sebastian Hunt, and Pasquale Malacaria. 2001. Quantitative Analysis of the Leakage of Confidential Data. *Electronic Notes in Theoretical Computer Science* 59, 3 (2001), 238–251.
- [13] David Clark, Sebastian Hunt, and Pasquale Malacaria. 2007. A static analysis for quantifying information flow in a simple imperative language. *Journal of Computer Security* 15, 3 (2007), 321–371.
- [14] Michael R. Clarkson and Fred B. Schneider. 2010. Hyperproperties. *Journal of Computer Security* 18, 6 (2010), 1157–1210.
- [15] Jared L. Cohon. 2004. *Multiobjective programming and planning*. Vol. 140. Courier Corporation.
- [16] Thomas M. Cover and Joy A. Thomas. 2006. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience.
- [17] Kalyanmoy Deb and Himanshu Jain. 2013. 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* 18, 4 (2013), 577–601.
- [18] Dorothy E. Denning. 1976. A Lattice Model of Secure Information Flow. *Commun. ACM* 19, 5 (1976), 236–243.
- [19] Andrea Fioraldi, Dominik Maier, Heiko Eißfeldt, and Marc Heuse. 2020. AFL++: Combining Incremental Steps of Fuzzing Research. In *14th USENIX Workshop on Offensive Technologies (WOOT 20)*. USENIX Association.
- [20] Gordon Fraser and Andrea Arcuri. 2011. 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*. 416–419.
- [21] Joseph A. Goguen and José Meseguer. 1982. Security Policies and Security Models. In *1982 IEEE Symposium on Security and Privacy*. 11–20.
- [22] S. He, M. Emmi, and G. Ciocarlie. 2020. ct-fuzz: Fuzzing for Timing Leaks. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*. IEEE Computer Society.
- [23] Jonathan Heusser and Pasquale Malacaria. 2010. Quantifying information leaks in software. In *Twenty-Sixth Annual Computer Security Applications Conference, ACSAC*. 261–269.
- [24] Pieter Hooimeijer, Martino Luca, Peter Oâ€™Hearn, Irene Papakonstantinou, Jim Purbrick, and Dulma Rodriguez. 2015. Moving Fast with Software Verification. In *NASA Formal Methods: 7th International Symposium, Proceedings*, Vol. 9058. Springer, 3.
- [25] Johannes Kinder. 2015. Hypertesting: The Case for Automated Testing of Hyperproperties. In *Hot Issues in security and trust (HotSpot)*.
- [26] Miquing Li, Tao Chen, and Xin Yao. 2020. How to Evaluate Solutions in Pareto-based Search-Based Software Engineering? A Critical Review and Methodological Guidance. *IEEE Transactions on Software Engineering* 01 (2020), 1–1.
- [27] LLVM Foundation 2022. Clang 15.0.0git documentation, AddressSanitizer. <https://clang.llvm.org/docs/AddressSanitizer.html>. Accessed: 2022-05-02.
- [28] Pasquale Malacaria, Michael Tautchning, and Dino Distefano. 2016. Information Leakage Analysis of Complex C Code and Its application to OpenSSL. In *Leveraging Applications of Formal Methods, Verification and Validation: Foundational Techniques - 7th International Symposium, ISOA 2016, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 9952)*. 909–925.
- [29] Valentin JM Manès, HyungSeok Han, Choongwoo Han, Sang Kil Cha, Manuel Egele, Edward J. Schwartz, and Maverick Woo. 2019. The art, science, and engineering of fuzzing: A survey. *IEEE Transactions on Software Engineering* 47, 11 (2019), 2312–2331.
- [30] Sergey Mehtaev, Jooyong Yi, and Abhik Roychoudhury. 2016. Angelix: Scalable multiline program patch synthesis via symbolic analysis. In *Proceedings of the 38th International Conference on Software Engineering (ICSE)*. 691–701.
- [31] Hector D. Menendez and David Clark. 2021. Hashing fuzzing: introducing input diversity to improve crash detection. *IEEE Transactions on Software Engineering* (2021).
- [32] Ibrahim Mesecan, Daniel Blackwell, David Clark, Myra B. Cohen, and Justyna Petke. 2021. HyperGL: Automated Detection and Repair of Information Flow Leakage. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1358–1362.
- [33] Jonathan Metzman, László Szekeres, Laurent Simon, Read Sprabery, and Abhishek Arya. 2021. FuzzBench: an open fuzzer benchmarking platform and service. In *ESEC/FSE*. 1393–1403.
- [34] Hoang Duong Thien Nguyen, Dawei Qi, Abhik Roychoudhury, and Satish Chandra. 2013. Semfix: Program repair via semantic analysis. In *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 772–781.
- [35] Yannic Noller and Saeid Tizpaz-Niari. 2021. QFuzz: Quantitative Fuzzing for Side Channels. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2021)*. ACM.
- [36] Justyna Petke, Saemundur O. Haraldsson, Mark Harman, William B. Langdon, David R. White, and John R. Woodward. 2017. Genetic improvement of software: a comprehensive survey. *IEEE Transactions on Evolutionary Computation* 22, 3 (2017), 415–432.
- [37] Andrei Sabelfeld and Andrew C. Myers. 2003. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications* 21, 1 (2003), 5–19.
- [38] Caitlin Sadowski, Edward Aftandilian, Alex Eagle, Liam Miller-Cushon, and Ciera Jaspan. 2018. Lessons from building static analysis tools at google. *Commun. ACM* 61, 4 (2018), 58–66.
- [39] Vali Tawosi, Federica Sarro, Alessio Petrozziello, and Mark Harman. 2021. Multi-objective software effort estimation: A replication study. *IEEE Transactions on Software Engineering* (2021).
- [40] The MITRE Corporation 2022. CVE – Common Vulnerabilities and Exposures. <https://cve.mitre.org/>. Accessed: 2022-05-02.
- [41] The National Institute of Standards and Technology 2014. CVE-2014-0160 Heartbleed. <https://nvd.nist.gov/vuln/detail/CVE-2014-0160>. Accessed: 2022-05-02.
- [42] The OpenSSL 2021. bn - multiprecision integer arithmetics. <https://www.openssl.org/docs/man1.0.2/man3/bn.html>. Accessed: 2022-03-12.
- [43] The OpenSSL 2021. OpenSSL 1.1.1 versions download. <https://ftp.openssl.org/source/old/1.1.1/>. Accessed: 2022-03-18.
- [44] Universal Ctags Team 2021. Universal Ctags. <https://ctags.io/>. Last Accessed: 2022-05-02.
- [45] Dennis M. Volpano, Cynthia E. Irvine, and Geoffrey Smith. 1996. A Sound Type System for Secure Flow Analysis. *Journal of Computer Security* 4, 2/3 (1996), 167–188.
- [46] Shuang Wang and Jeff Offutt. 2009. Comparison of unit-level automated test generation tools. In *2009 International Conference on Software Testing, Verification, and Validation Workshops*. IEEE, 210–219.
- [47] David R. White, Andrea Arcuri, and John A. Clark. 2011. Evolutionary improvement of programs. *IEEE Transactions on Evolutionary Computation* 15, 4 (2011), 515–538.
- [48] Hirotooshi Yasuoka and Tachio Terauchi. 2014. Quantitative information flow as safety and liveness hyperproperties. *Theoretical Computer Science* 538 (2014), 167–182.