

Leveraging Program Invariants to Promote Population Diversity in Search-Based Automatic Program Repair

Zhen Yu Ding
School of Computing and Information
University of Pittsburgh
Pittsburgh, PA, USA
zhd23@pitt.edu

Yiwei Lyu, Christopher S. Timperley and Claire Le Goues
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, USA
ylyu1@andrew.cmu.edu, {ctimperley, clegoues}@cs.cmu.edu

Abstract—Search-based automatic program repair has shown promise in reducing the cost of defects in real-world software. However, to date, such techniques have typically been most successful when constructing short or single-edit repairs. This is true even when techniques make use of heuristic search strategies, like genetic programming, that in principle support the construction of patches of arbitrary length. One key reason is that the fitness function traditionally depends entirely on test cases, which are poor at identifying partially correct solutions and lead to a fitness landscape with many plateaus. We propose a novel fitness function that optimizes for both functionality and semantic diversity, characterized using learned invariants over intermediate behavior. Our early results show that this new approach improves semantic diversity and fitness granularity, but does not statistically significantly improve repair performance.

Index Terms—genetic programming, automated program repair, program semantics, invariant analysis, fitness function

I. INTRODUCTION

Research in automated program repair aims to reduce the significant cost of buggy software in terms of both societal impact [1] and developer effort [2]. Most modern repair techniques take, as input, source code and a test suite, and seek a patch that leads all input tests to pass, including tests that initially fail. One dominant class of repair approaches relies on heuristic search strategies like random search [3], a deterministic heuristic walk of the search repair space [4]–[6], or genetic programming [7], [8]. These techniques are promising in their application to both historical bugs in real-world programs as well as, more recently, industrial settings [9]. However, to date, most techniques are limited to finding single-edit repairs. For some techniques, this is by design [3], [4], [8], and avoids the challenges of traversing a combinatorial multi-edit search space. However, even techniques like GenProg, which can in principle construct multi-edit repairs, still often produce patches that reduce to single edits [10]. This inability to construct multi-edit repairs limits the scope of such repair techniques to simpler, single-line bugs.

One major factor that limits the ability of genetic programming to evolve multi-edit bug repairs is a fitness function that relies on test cases. Ideally, a candidate patch that contains

a partial solution—such as one of several required edits—would cause the program to pass more test cases than one that contains no useful edits. The test-based fitness function could then usefully identify partial solutions, which could then be selected and recombined for evolution towards a full solution. Unfortunately, this does not typically hold in practice. Instead, test case-based fitness functions often assign the same fitness score to many different candidate solutions, producing large *fitness plateaus* [11], [12]. An insufficiently granular fitness function that fails to distinguish between distinct candidate solutions prevents the search from identifying promising edits, and instead simply prevents the introduction of destructive edits that cause test failures [13]. Moreover, tests inherently lack information about intermediate program behavior. They are therefore poorly suited to guiding a search towards semantically diverse mechanisms to accomplishing the same repair, a possibly fruitful application for program repair [14].

We seek to overcome these limitations by providing additional information about program semantics to help the genetic programming search strategy distinguish between promising candidate patches. One potential source of additional information is *likely program invariants* over program values, like $\mathbf{x} == \mathbf{0}$ or $\mathbf{y} > \mathbf{a}$. There exist mature dynamic invariant detection tools, such as Daikon [15], that infer likely invariants observed as holding over multiple program executions.

We propose to infer a set of likely invariants by observing the *original* program behavior on its *initially passing* test cases, which demonstrate good behavior. We then evaluate the behavior of the program as modified by candidate patch solutions on those inferred invariants. We use this observed behavior over the population of candidate patches to determine which patch variants appear to be more *semantically diverse*. We then use semantic diversity as an optimization objective, along with the number of passing test cases, to determine variant fitness.

Beyond the potential utility of a technique that can find diverse solutions to a program repair problem, diversity among variants in a genetic programming population discourages genetic drift and premature convergence towards local op-

tima [16]. By injecting semantic diversity, we hope that our new fitness function can more effectively rank candidate patches that would otherwise appear identical to a purely test case-based fitness function, leading to higher fitness granularity. Our goal is a repair search strategy with fewer fitness plateaus and can generate semantic diversity in its population of candidate solutions. Our technique is designed to be more general than previous attempts to either examine memory values [12], or to use program predicates [17] to characterize repair fitness, and further does not require training data.

This paper’s main contributions are:

- A new repair technique that uses learned program invariants to promote semantic diversity of variant programs.
- A compact, invariant-based representation of program semantics. Our representation facilitates an approximation of the semantic difference between programs.
- An evaluation of an implementation of our technique’s performance across sampled IntroClassJava [18] bugs.
- An evaluation of our technique’s scalability to Defects4j [19], a set of larger, real-world buggy programs.

The rest of this paper proceeds as follows. Section II outlines background helpful to understanding our approach. We motivate the approach with an illustrative example in Section III and then describe the approach in Section IV. Section V describes our evaluation on both large and small programs, showing our technique’s potential to reduce fitness plateaus and increase population diversity. We describe threats to the validity of our experiments and possible limitations of our approach in Section VI. Section VII puts our work in context with respect to the prior literature; we conclude with a discussion and vision for future directions in Section VIII.

II. BACKGROUND

This section provides background helpful to understanding our technique and contributions.

a) Automatic Program Repair (APR): Given a program and an oracle that determines its correctness, an Automatic Program Repair technique attempts to find a *patch* that modifies the program in a way the oracle deems correct. The majority of modern repair techniques use test cases as an oracle, as a widely available but incomplete proxy for formal correctness specifications. Test cases that fail on the original program (*negative* tests) expose the bug to be repaired; test cases that pass on the original program (*positive* tests) indicate desirable behavior that should be maintained. Thus, the goal of APR is to find a patch such that all provided tests pass on the modified program. We describe the previous alternatives to this approach in Section VII. Repair techniques typically also use the test cases to perform fault localization (generally using an off-the-shelf technique [20]), constraining the search to a smaller set of candidate repair locations.

Two dominant approaches to APR are *semantics-based* and *heuristic* or *search-based*. Semantics-based tools [21]–[23] seek to synthesize a repair by deducing intended program semantics. Heuristic or search-based approaches explore search spaces of templated repairs applied to the abstract syntax tree

of the program. Heuristic techniques [4]–[6] use predefined schemas and probabilistic models or other heuristics to apply candidate changes. We focus in this work on search-based techniques like GenProg [7] and RSRepair [3] that explicitly use stochastic techniques to traverse the search space.

Techniques in both classes have been successfully applied to real-world programs. However, they vary in the type of code they can handle and the types of patches they can create. Semantics-based approaches are often restricted by their underlying reasoning to synthesizing new conditionals or, less commonly, the right-hand side of assignments. Search-based techniques can provide a broader set of syntactic repair templates, but must still be restricted in the type and variety that they explore, or the search becomes intractable. Additionally, many techniques are expressly limited to single-edit repairs. Even techniques (like GenProg, Angelix, or S3) that can in principle construct multi-edit repairs rarely do, and seemingly multi-edit patches often reduce to a single edit [10].

b) GenProg: GenProg [7] is a search-based APR tool that uses genetic programming [24] to traverse the space of possible patches. Patches consist of a variable-length sequence of *statement-level* edits to the original program (typically *deletion*, *replacement*, or *insertion*). Edits requiring new code (insertion or replacement) restrict the domain of new code to elsewhere in the same buggy program. New patches are created via recombination (using *crossover*) and mutation (by appending a new randomly-instantiated edit to an existing patch candidate). Intermediate solutions are *selected* to continue evolution to subsequent generations, typically with a selection scheme weighted by variant *fitness*, or suitability. Like most other repair techniques, GenProg optimizes for patched program performance on the input test cases, with a fitness function comprised of a weighted sum of the number of positive and negative tests that a modified program passes. Empirically, however, this fitness function lacks granularity and often fails to precisely distinguish between partial solutions [11], [12], [17]. GenProg was initially implemented to repair C programs, but has been reimplemented for Java.¹

c) Automatic Learning of Likely Program Invariants: There are tools that automatically infer likely program invariants, or predicates that describe program behavior. Such techniques have been used to characterize either normal or abnormal program executions, useful in applications like specification mining [15], N-variant systems [25], and fault localization [26]. We use Daikon [15], a mature dynamic analysis technique that infers likely program invariants by running the program on a provided set of inputs or tests, observing intermediate values of the program, and matching those values to template predicates to report properties holding true over all observed executions. Daikon’s released toolset² works on a variety of languages, including Java. We discuss Daikon’s limitations in Section VI.

```

1 public int gcd(int a, int b) {
2     int result = 1;
3     if (a == 0) {
4         b = b - a;
5     } else {
6         result=a;
7         while (b != 0) {
8             result = b;
9             if (a > b) {
10                a = a - b;
11            } else {
12                b = b - a;
13            }
14        }
15    }
16    result=a;
17    return result;
18 }

```

Fig. 1: A buggy implementation of Euclid’s GCD algorithm. This function is incorrect when $a = 0$ (returns 0, not b).

III. MOTIVATING EXAMPLE

We begin with a short example to motivate our proposed approach. Consider the Java function implementing Euclid’s GCD algorithm shown in Figure 1. This code contains a two-part bug: when $a = 0$, it will return 0 instead of b (the correct answer); when $b = 0$ and $a > 0$, it will return a (instead of 0). A simple patch for this bug requires two changes: deletion of line 16, and a replacement of line 4 with line 8 (or appending line 8 before or after line 4). Note that GenProg can in principle construct this patch using statement-level modifications: GenProg’s mutation operators can for example delete all of line 16 in Figure 1, or replace the statement at line 4 with the one at line 8 (but could not modify the expression in any of the if conditions on lines 3 or 9).

Although test cases usefully characterize program behavior, they typically do not offer information on partial program correctness. More importantly, two very different programs can easily pass the same number of test cases. For example, a candidate patch that simply deletes line 16 would result in a program that passes the same tests as the original buggy code. However, this deletion is one of the two required edits. Ideally, the fitness function would rank it above the empty patch.

Program invariants offer additional semantic information that could distinguish between such patches. Programs that display different behavior with respect to a set of invariants may be interestingly different from one another in terms of intermediate program semantics. For example, given a small set of positive test cases, Daikon can infer likely invariants such as $a \% \text{gcd}(a, b) = 0$ and $b \% \text{gcd}(a, b) = 0$. These invariants describe important properties of the `gcd` function.

Note that the buggy code in Figure 1 violates these properties on an initially failing test where $a = 0$. On the other hand,

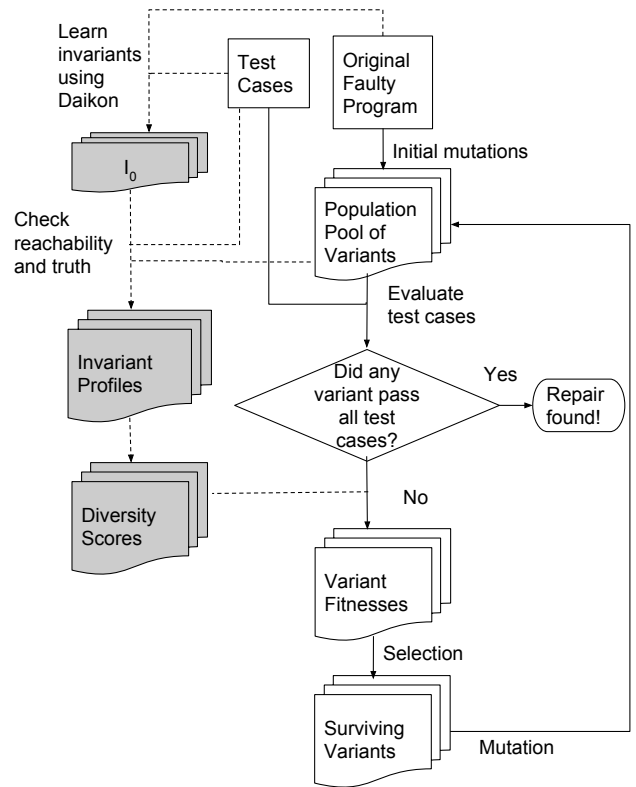


Fig. 2: A flow-chart of our approach: White boxes and solid lines show GenProg’s architecture, grey boxes and dotted lines are our approach’s additions.

the candidate patch that deletes line 16 produces a program that correctly preserves those invariants. This suggests that the change may be worth maintaining in the search.

Note that we do not know a priori which possible properties are most important, and Daikon can produce noisy output with trivial or overfitting properties. Fortunately, genetic programming [27], [28] is robust to noisy fitness functions. Thus, we propose in this initial work to use inferred likely invariants as one component of a genetic programming-based repair search to promote diversity and reduce fitness plateaus.

IV. APPROACH

In this section, we describe our new genetic programming-based repair technique that optimizes for both *functionality* (measured using test cases) and *semantic diversity*. We measure diversity by using inferred program invariants to characterize modified program (and thus candidate patch) behavior. We call this its *invariant profile* (Section IV-A). We then quantify a candidate patch’s diversity by comparing its invariant profile to the rest of the patches in the search population (Section IV-B). We use this invariant-based diversity score, along with the number of positive and negative test cases passed, as optimization objectives in the multi-objective genetic algorithm NSGA-II [29] (Section IV-C).

Figure 2 shows the high-level architecture of our approach. We add the following to GenProg’s core approach: First, we

¹<https://github.com/squaresLab/genprog4java>

²<https://plse.cs.washington.edu/daikon/>

run Daikon on the original faulty program with positive test cases to obtain a set of inferred possible invariants, denoted by I_0 . Then, in each generation, we run tests and instrumented code through each variant in the population to obtain invariant profiles. We calculate diversity scores for each variant in the population based on this profile, and use the diversity scores in the fitness function in the multi-objective search.

A. Invariant Profiles

For each candidate patch, we construct an invariant profile containing information about the patch’s behavior. This profile consists of the reachability and truth value of a set of potentially interesting invariants inferred by Daikon [15]. Due to Daikon’s high computational cost, we do not learn invariants from each patch. Instead, we infer potentially interesting invariants from the original buggy program to create a starting set of invariants. Moreover, we only run Daikon on positive tests to construct this set of invariants, due to the practical limitations of running Daikon on failing tests that crash the program. We believe that learning properties over negative tests may also be interesting, and leave the engineering and conceptual approach to doing so to future work.

Daikon can learn potential invariants at multiple levels of program abstraction. However, because patches can manipulate code arbitrarily, predicates that are sensible in describing original intermediate program behavior may no longer be evaluable on patched code. We therefore include only *method-level* invariants in the starting set, denoted by I_0 , which enables evaluation of whether an invariant continues to hold in an arbitrary candidate patch.

Even though we use only positive tests for learning the starting set of invariants, we do check whether invariants in I_0 hold over both positive and negative tests when constructing invariant profiles for candidate patches.

Assuming the following definitions:

$$\begin{aligned} T &= \text{Set of all test cases in a program’s test suite} \\ T^+ &= \text{Set of all positive test cases in } T \\ T^- &= \text{Set of all negative test cases in } T \\ I_0 &= \text{Set of inferred method-level invariants} \end{aligned}$$

The invariant profile for a particular patch v is constructed as follows. We apply v to the original program code to produce a variant program. We then instrument the variant program’s code to evaluate, for every invariant p in I_0 , whether the method associated with p is ever reached, and if so, whether p holds for the variant or not. We collect the aforementioned data separately for T^+ and T^- , and track and incorporate the invariant data derived from the two sets of tests separately.

Let p_i represent the i th invariant in I_0 . We represent the invariant profile of a variant m , Π_m , as a ternary vector of length $2|I_0|$, where for all i between 0 and $|I_0|$, the $(2i)$ -th and $(2i + 1)$ -th characters of Π_m are:

$$\Pi_m[2i] = \begin{cases} 0 & p_i \text{ is always true for all tests in } T^+ \\ 1 & p_i \text{ is false at least once in } T^+ \\ 2 & p_i \text{ is never reached by any test in } T^+ \end{cases}$$

$$\Pi_m[2i + 1] = \begin{cases} 0 & p_i \text{ is always true for all tests in } T^- \\ 1 & p_i \text{ is false at least once in } T^- \\ 2 & p_i \text{ is never reached by any test in } T^- \end{cases}$$

B. Quantifying Invariant-based Diversity

To score the invariant-based diversity of candidate patches, we compare a patch’s invariant profile to that of every other patch in the population. The goal is to determine the uniqueness of a patch’s invariant profile with respect to the other candidate patches. To measure the difference in invariant behavior between two patches m_1 and m_2 , we use the Hamming distance of their invariant profiles Π_{m_1} and Π_{m_2} :

$$\Delta(m_1, m_2) = \text{HammingDistance}(\Pi_{m_1}, \Pi_{m_2})$$

The diversity score of a patch m in population P , assuming that m compiles and terminates within a specified time limit, is the sum of the Hamming distances between m ’s invariant profile and the profiles of every other patch in P whose variant programs compile and terminate:

$$\text{diversity}(m) = \sum_{n \in P} \begin{cases} \Delta(m, n) & m \text{ and } n \text{ compile and terminate} \\ 0 & \text{otherwise} \end{cases}$$

In general, diversity among variants in a population discourages premature convergence towards local optima [16]. By promoting semantic diversity, we aim to encourage a repair-oriented genetic algorithm to explore niches in semantic space.

C. Promoting Diversity

We use the invariant-based semantic diversity score as an optimization objective, alongside the number of positive and negative tests passed. This treats program repair as a multi-objective optimization problem. We use the NSGA-II algorithm [29] for this purpose. NSGA-II is a genetic algorithm that optimizes for multiple objectives using Pareto optimality. NSGA-II also employs mechanisms to encourage population diversity over the course of its search.

Note that we optimize for positive and negative tests as separate objectives rather than using a weighted sum of positive/negative tests as a single objective, as GenProg does.

V. EVALUATION

To analyze the effects of promoting invariant-based diversity, we implemented our proposed technique as an extension of GenProg4J³ and used the original version of GenProg4J as a control. Using our modified version of GenProg, we answer the following questions:

- RQ1** Does optimizing for invariant-based diversity change the effectiveness in searching for a repair with respect to the number of repairs found, search performance, and search success rate?
- RQ2** Does optimizing for invariant-based diversity increase the diversity of patches in a population?
- RQ3** Does invariant-based diversity provide more granular fitness (and therefore less plateauing)?

³<https://github.com/squaresLab/genprog4java>

Bug Scenario		GenProg	Our Approach
program	ID		
smallest	af81-000	10%	10%
	8839-002	100%	100%
	0cdf-006	100%	100%
digits	0cdf-007	90%	100%
	d120-001	80%	90%
	5b50-000	10%	0%
Total		65%	67%

TABLE I: Repair success rates for each bug repaired by either GenProg or our approach. The unlisted 53 bugs have a success rate of 0% for both techniques.

RQ4 Does our diversity-enhanced extension of GenProg scale to large, real-world programs?

A. Efficiency, Repairability and Diversity

To answer RQs 1, 2, and 3, we use a modified version of the IntroClassJava [18] dataset of bugs (see Section VI for a discussion of the changes). IntroClassJava is a set of small Java programs (< 30 lines) derived from IntroClass [30], a set of buggy C programs written by students in an introductory CS course. We use IntroClassJava as a benchmark since its design and the small size of its constituent programs supports controlled experimentation via larger numbers of repairs.

We ran GenProg and our technique on a random sample of 59 out of a possible 297 bugs from IntroClassJava, and repeated each experiment with 10 different random seeds. To ensure a fair comparison between GenProg and our technique, we use an identical configuration for both techniques, based on the default settings for GenProg4J. We used the append, replace, and delete repair operators, one-point crossover, tournament selection ($k = 8$),⁴ a population size of 40, and ran for 10 generations. We ran each experiment on Amazon EC2 using a *c5.l* instance with two vCPUs, 4 GB of RAM, and a 100 GB *gp2* storage volume, running Ubuntu 16.04 LTS.

1) *RQ1 – Performance*: Across both GenProg and our approach, a total of six unique bugs were successfully repaired at least once. Table I provides the repair success rate (i.e., the likelihood that a given attempt to repair the bug will be successful) for each of the six bugs.

We find that, with one exception, our approach and GenProg fix the same bugs. In the exceptional case, GenProg repairs `digits-5b50-000` once in ten seeds. Using a two-sided version of Fisher’s exact test, we fail to demonstrate a statistically significant difference ($p < 0.05$) in the success rate between GenProg and our approach for each individual bug and for all bugs in aggregate. Given GenProg’s low success rate, our technique’s inability to find a repair on `digits-5b50-000` may be due to a sampling error.

Figure 3 compares the *efficiency* of our approach against GenProg for the five bugs that were repaired by both techniques, where we define efficiency as the number of unique

⁴At the time of our experiments, the default tournament size in GenProg4J is 20% of the population size.

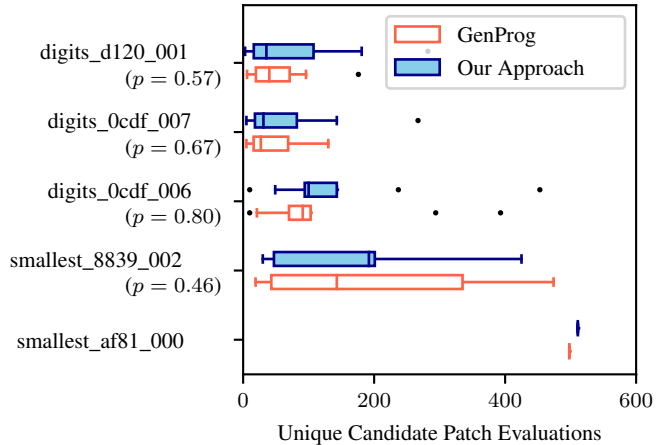


Fig. 3: A comparison of the number of unique candidate patch evaluations before a repair was found for the five bugs fixed by both our approach and GenProg. Using a one-sided Mann-Whitney U test, we fail to demonstrate a statistically significant difference ($p < 0.05$) for any bug.

Bug Scenario		GenProg	Our Approach
program	ID		
smallest	af81-000	1 / 1	1 / 1
	8839-002	8 / 10	8 / 10
	0cdf-006	9 / 10	10 / 10
digits	0cdf-007	8 / 9	9 / 10
	d120-001	4 / 8	6 / 9
	5b50-000	1 / 1	0 / 0
Total		31 / 39	34 / 40

TABLE II: Number of unique repairs found for each bug across all seeds.

candidate patch evaluations. Using a one-sided Mann-Whitney U test [31], we fail to find a statistically significant improvement ($p < 0.05$) over GenProg across all bugs.

Table II reports the number of syntactically unique patches for each bug. Using a two-sided version of Fisher’s exact test, we were unable to demonstrate a statistically significant difference at the $\alpha = 0.05$ confidence level ($p = 0.249$).

In summary, our results fail to demonstrate a significant difference, be it negative or positive, to efficiency, repair success rate, and the number of unique patches.

2) *RQ2 – Diversity*: To determine whether our search technique is an effective means of promoting semantic diversity within the population, we compare the *normalized population diversity* (NPD) between GenProg and our approach across all 59 bugs sampled from IntroClassJava. We define the NPD of a given population P as the sum of diversity scores of the individual candidate patches belonging to that population, divided by the number of individual patches in the population squared and the length of invariant profile (Π) for a given bug.

$$NPD(P) = \frac{\sum_{m \in P} \text{diversity}(m)}{|P|^2 \cdot |\Pi|}$$

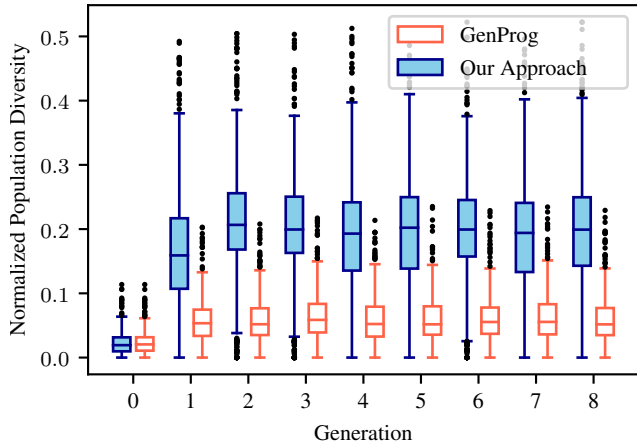


Fig. 4: A comparison of normalized population diversity. Except in generation 0, there is a statistically significant difference ($p < 0.001$), computed using a one-sided Mann-Whitney U test, between GenProg and our approach.

Figure 4 shows the distribution of NPD scores for GenProg and our approach at each generation prior to selection. Our results show significantly higher diversity within the population (after the initial population) with our approach, suggesting that our algorithm is effective at promoting diversity. We observe low levels of population diversity in the initial population of the search (i.e., generation 0) for both our approach and GenProg; this result is expected as the population has not been subject to selection. With our approach, diversity levels increase until the second generation of the search, at which point diversity stabilizes.

For GenProg, we observe a significant increase in diversity between generation 0 and all subsequent generations, which immediately stabilizes after the first generation. This result suggests that GenProg’s test-based fitness function implicitly selects for semantic diversity, albeit relatively weakly.

3) *RQ3 – Fitness Granularity*: We measure the granularity of fitness information available to each repair method by analyzing the phenotypic diversity [32], [33] of populations. We use a ratio of unique fitness scores to the size of the population as a measure of phenotypic diversity.

Figure 5 compares fitness granularity between GenProg and our approach. Using a one-sided Mann-Whitney U test, we find that our approach significantly ($p < 0.001$) increases fitness granularity. We observe a 71% increase in median fitness granularity from 0.088 to 0.150 with our approach.

B. Scalability

To answer RQ4, we compare the wall-clock time for GenProg and our approach to run on 10 bugs from the Defects4J [19] dataset of real-world Java bugs. To accommodate the larger test suites of Defects4J programs, we no longer infer invariants I_0 from the full set of positive tests T^+ . Instead, we use a subset of T^+ that only includes positive tests that co-occur in JUnit test classes that also contain negative tests.

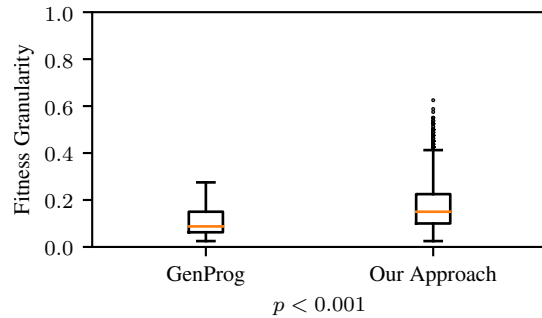


Fig. 5: A comparison of fitness granularity within the population prior to selection in each generation.

Bug	GenProg	Our Approach	Difference
lang11	3586 s	3862 s	1.08 X
lang29	1783 s	2223 s	1.25 X
lang36	2088 s	2465 s	1.18 X
lang8	5850 s	6239 s	1.07 X
lang9	3304 s	4252 s	1.29 X
math30	5356 s	5433 s	1.01 X
math44	5906 s	10613 s	1.80 X
math46	4023 s	43229 s	10.75 X
math79	6033 s	7178 s	1.19 X
math86	3751 s	4287 s	1.14 X
Median	3887 s	4860 s	1.19 X
Mean	4168 s	8978 s	2.18 X

TABLE III: A comparison of the wall-clock time GenProg and our approach takes to run a single seed across 10 bugs from the Defects4J benchmark, in seconds.

We use a near-identical configuration to our IntroClassJava experiments with two changes: (1) we use a population size of 20, and (2) we evaluate variants on a random sample of 10% of the positive tests [17], along with all negative tests. We ran all experiments on a machine with the following specifications: Ubuntu 16.04 LTS, 4 x Intel Xeon E7-4820 (40 cores), 128 GB RAM, and 2 TB storage.

Table III presents a summary of our results. We find that the median run-time is 19% higher when using our approach. We expect our approach to be naturally slower than GenProg since it introduces two notable overheads: (1) before the search, we run Daikon on the original program to infer a set of likely invariants I_0 ; (2) we instrument each variant program and re-run the test suite to obtain an invariant profile. The static overhead from (1) is necessary only once per bug and amortizable across multiple repair attempts. We can reduce the dynamic overhead from (2) by further optimizing our implementation (e.g., by avoiding redundant test executions).

The order-of-magnitude increase in run-time for `math46` is likely explained by an abnormally large number of inferred invariants. The object file containing all invariants instrumented for `math46` is 1.91MB, an order of magnitude larger than the object files of other `math` bugs’ instrumented invariants (varies between 16KB and 232KB). Our implementation adds instrumentation for each individual predicate in I_0 to the

program. As a result, a large amount of instrumentation was added to `math46`, creating a substantial static overhead. We may reduce this overhead by optimizing the instrumentation process and by filtering spurious invariants.

In summary, our results show that our prototype implementation, in most cases, introduces a tractable run-time computational overhead (approximately 19%) over base GenProg, suggesting it may plausibly be applied to larger programs.

VI. THREATS AND LIMITATIONS

a) Manual Code Modification: In porting IntroClass for C [30] to Java, the IntroClassJava authors elected to include all primary functionality in a single Java method. Input values are retrieved using I/O in the middle of this single method. Daikon struggles to identify interesting invariants over code organized this way. We thus mechanically refactored IntroClassJava code to relocate core logic into a separate method and more clearly delineate input and output values for the different problems. This decouples our technique performance from Daikon performance. We anticipate that dataflow analysis or slicing techniques can automate the task of identifying potentially interesting input and output values, as has been done in different program repair contexts [34].

b) Repair Quality: Repair quality is a key unsolved concern in program repair. Both search- and semantics-based approaches can create patches that *overfit* to the provided test cases, or fail to generalize to the desired but unwritten program specification [35], [36]. We did not evaluate the quality of the repairs that GenProg and our repair technique produced, relying instead on performance on the provided test cases. We leave assessment of repair quality to future work.

c) Invariant Quality: While Daikon can detect various types of invariants, Daikon’s limited set of supported invariants and scope restrict the tool’s flexibility. Daikon’s inability to detect invariants relating to the state of variables at intermediate statements of a method lead us to refactor IntroClassJava.

Daikon may also produce trivial invariants such as `this != null`. These invariants provide no useful semantic information. Unnecessarily evaluating these invariants wastes computational resources. Inferred invariants may also overfit. For example, if a program executes the function `sin(x)` only ever with input `x = pi * n`, Daikon might infer `sin(x) = 0` as an invariant. Nonetheless, overfit invariants may sometimes be useful in our approach, as breaking them with the same input indicates a change in program semantics. We may attempt to mitigate these limitations in the future by filtering Daikon’s output to remove low-quality invariants.

VII. RELATED WORK

There exist previous efforts to create better fitness functions using invariants [17], intermediate program state [12], and historical repair patterns [8] to guide the search.

Fast et al. [17] propose an alternative fitness function for GenProg that uses learned invariants to identify and promote partial solutions. Their fitness function is implemented as a linear model that accepts inferred invariants and test cases as

its features. Existing patches are used to create training data that is used to construct a linear model for a given program and to learn associations between invariants and patch correctness. In contrast to Fast et al.’s approach, our approach does not require training data.

de Souza et al. [12] aimed to reduce fitness plateaus by injecting source code checkpoints to track pre-determined variables. Based on changes in these variables and the results of fault localization, they incorporate a *checkpoints* metric into the fitness function. de Souza et al. found that their approach reduced plateauing, fixed more bugs, and improved efficiency. Our approach provides an orthogonal method for reducing plateauing that uses automatically learned invariants instead of the values of pre-determined variables.

Timperley et al. [37] aim to reuse variant test results collected over the course of the search to determine fix locations via a dynamic mutation-based fault localization. They found that modifications to correct statements are more likely to result in the failure of previously passing tests, but that few previously passing tests covered the faulty area of the program. Ultimately, they were unable to demonstrate a significant improvement to the accuracy of fault localization when variant test cases were analyzed. Our approach might address the lack of sufficient test coverage necessary to perform effective mutation analysis — at no additional cost — by observing changes to invariant profiles rather than test cases.

Beadle and Johnson [38] investigated promoting semantic diversity during population initialization in genetic programming. They developed semantically driven initialization algorithms for Boolean and artificial ant problems. Using their algorithms often resulted in better performance compared to RHH [24], a standard initialization technique. Their results suggest that promoting invariant diversity during initialization may yield greater performance.

VIII. CONCLUSIONS

Our early results demonstrate that inferred invariants can be used to effectively increase semantic diversity within the population, and to reduce fitness plateauing in search-based automatic program repair. Using Defects4J, we demonstrate that a prototype implementation of our technique scales to bugs in large-scale, real-world programs. In our experiments on IntroClassJava, we fail to demonstrate a statistically significant improvement to efficiency, the number of bugs fixed, unique patches found, and repair success rate.

A possible explanation for this result is that the search space contains few, if any, acceptable repairs. Our approach uses the same statement-level repair operators as GenProg, and so it shares the limitations of GenProg’s search space. To obtain a better idea of how our approach improves search performance, we plan to evaluate our approach with a richer search space, generated using a larger number of mutation operators.

Going forward, we also plan to collect more data on our approach’s performance on a larger number of bugs, including the rest of the IntroClassJava and Defects4J datasets, along with JaConTeBe [39] and Bugs.jar [40]. We may also attempt

to improve our approach's run time by pruning invariants and by optimizing the code instrumentation process.

To encourage further investigation and ensure reproducibility, we provide a prototype implementation of our technique, along with the results and raw data from our study, as part of our replication package at <https://bit.ly/2HNQ5g7>

ACKNOWLEDGEMENTS

This work is partially supported by the NSF under grants CCF-170116, CCF-1833698, and REU Site CNS-1560137.

REFERENCES

- [1] "Software fail watch: 5th edition," <https://www.tricentis.com/software-fail-watch/>, Accessed December, 2018.
- [2] "Cambridge university study states software bugs cost economy \$312 billion per year," <http://www.prweb.com/releases/2013/1/prweb10298185.htm>, Accessed December, 2018.
- [3] Y. Qi, X. Mao, Y. Lei, Z. Dai, and C. Wang, "The strength of random search on automated program repair," in *International Conference on Software Engineering*, ser. ICSE '14, 2014, pp. 254–265.
- [4] W. Weimer, Z. P. Fry, and S. Forrest, "Leveraging program equivalence for adaptive program repair: Models and first results," in *International Conference on Automated Software Engineering*, ser. ASE '13, 2013, pp. 356–366.
- [5] F. Long and M. Rinard, "Staged program repair with condition synthesis," in *Joint Meeting on Foundations of Software Engineering*, ser. ESEC/FSE '15, 2015, pp. 166–178.
- [6] —, "Automatic patch generation by learning correct code," in *Principles of Programming Languages*, ser. POPL '16, 2016, pp. 298–312.
- [7] C. Le Goues, T. Nguyen, S. Forrest, and W. Weimer, "GenProg: A generic method for automatic software repair," *IEEE Trans. Software Eng.*, vol. 38, pp. 54–72, 2012.
- [8] X. D. Le, D. Lo, and C. L. Goues, "History driven program repair," in *International Conference on Software Analysis, Evolution, and Reengineering*, ser. SANER '16, vol. 1, 2016, pp. 213–224.
- [9] A. Marginean, J. Bader, S. Chandra, M. Harman, Y. Jia, K. Mao, A. Mols, and A. Scott, "SapFix: Automated end-to-end repair at scale," in *International Conference on Software Engineering (SEIP track)*, ser. ICSE-SEIP '19, 2019, p. (to appear).
- [10] Z. Qi, F. Long, S. Achour, and M. Rinard, "An analysis of patch plausibility and correctness for generate-and-validate patch generation systems," in *International Symposium on Software Testing and Analysis*, ser. ISSTA '15, 2015, pp. 24–36.
- [11] S. Forrest, W. Weimer, T. Nguyen, and C. Le Goues, "A genetic programming approach to automated software repair," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2009, pp. 947–954.
- [12] E. F. de Souza, C. Le Goues, and C. G. Camilo-Junior, "A novel fitness function for automated program repair based on source code checkpoints," in *Genetic and Evolutionary Computation Conference*, ser. GECCO '18, 2018.
- [13] C. S. Timperley, "Advanced techniques for search-based program repair," Ph.D. dissertation, University of York, June 2017.
- [14] R. Feldt, "Generating diverse software versions with genetic programming: An experimental study," *IEEE Proceedings-Software*, vol. 145, no. 6, pp. 228–236, 1998.
- [15] M. D. Ernst, J. H. Perkins, P. J. Guo, S. McCamant, C. Pacheco, M. S. Tschantz, and C. Xiao, "The Daikon system for dynamic detection of likely invariants," *Science of Computer Programming*, vol. 69, no. 1–3, pp. 35–45, 2007.
- [16] A. E. Eiben and J. E. Smith, *Introduction to evolutionary computing*. Springer, 2011.
- [17] E. Fast, C. Le Goues, S. Forrest, and W. Weimer, "Designing better fitness functions for automated program repair," in *Genetic and Evolutionary Computation Conference*, ser. GECCO '10, 2010, pp. 965–972.
- [18] T. Durieux and M. Monperrus, "IntroClassJava: A benchmark of 297 small and buggy Java programs," Universite Lille 1, Research Report, 2016.
- [19] R. Just, D. Jalali, and M. D. Ernst, "Defects4J: A database of existing faults to enable controlled testing studies for Java programs," in *International Symposium on Software Testing and Analysis (ISSTA)*, 2014, pp. 437–440.
- [20] J. A. Jones and M. J. Harrold, "Empirical evaluation of the Tarantula automatic fault-localization technique," in *Automated Software Engineering*, 2005, pp. 273–282.
- [21] J. Xuan, M. Martinez, F. Demarco, M. Clement, S. R. L. Marcote, T. Durieux, D. Le Berre, and M. Monperrus, "Nopol: Automatic repair of conditional statement bugs in Java programs," *IEEE Trans. Software Eng.*, vol. 43, no. 1, pp. 34–55, 2017.
- [22] X. D. Le, D. H. Chu, D. Lo, C. Le Goues, and W. Visser, "S3: Syntax- and semantic-guided repair synthesis via programming by examples," in *Joint Meeting of the European Software Engineering Conference and the Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE '17, 2017, pp. 593–604.
- [23] S. Mechtaev, J. Yi, and A. Roychoudhury, "Angelix: Scalable multiline program patch synthesis via symbolic analysis," in *International Conference on Software Engineering*, ser. ICSE '16, 2016, pp. 691–701.
- [24] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, 1992.
- [25] J. H. Perkins, S. Kim, S. Larsen, S. Amarasinghe, J. Bachrach, M. Carbin, C. Pacheco, F. Sherwood, S. Sidiroglou, G. Sullivan, W. Wong, Y. Zibin, M. D. Ernst, and M. Rinard, "Automatically patching errors in deployed software," in *Symposium on Operating Systems Principles*, 2009, pp. 87–102.
- [26] B. Liblit, M. Naik, A. X. Zheng, A. Aiken, and M. I. Jordan, "Scalable statistical bug isolation," in *Conference on Programming Language Design and Implementation*, ser. PLDI '05, 2005, pp. 15–26.
- [27] J. R. Koza, "Genetic programming as a means for programming computers by natural selection," *Statistics and Computing*, vol. 4, no. 2, pp. 87–112, Jun 1994.
- [28] B. L. Miller and D. E. Goldberg, "Optimal sampling for genetic algorithms," University of Illinois at Urbana-Champaign, Tech. Rep., August 1996.
- [29] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: NSGA-II," *IEEE Trans. Evolutionary Comp.*, vol. 6, no. 2, pp. 182–197, 2002.
- [30] C. Le Goues, N. Holtschulte, E. K. Smith, Y. Brun, P. Devanbu, S. Forrest, and W. Weimer, "The ManyBugs and IntroClass benchmarks for automated repair of C programs," *IEEE Trans. Software Eng.*, vol. 41, no. 12, pp. 1236–1256, 2015.
- [31] H. B. Mann and D. R. Whitney, "On a test of whether one of two random variables is stochastically larger than the other," *The Annals of Mathematical Statistics*, vol. 18, no. 1, pp. 50–60, 1947.
- [32] P. D'haeseleer, "Context preserving crossover in genetic programming," in *Conference on Evolutionary Computation. World Congress on Computational Intelligence*, vol. 1, 1994, pp. 256–261.
- [33] E. Burke, S. Gustafson, and G. Kendall, "A survey and analysis of diversity measures in genetic programming," in *Genetic and Evolutionary Computation Conference*, ser. GECCO'02, 2002, pp. 716–723.
- [34] Y. Ke, K. T. Stolee, C. Le Goues, and Y. Brun, "Repairing programs with semantic code search," in *Automated Software Engineering*, ser. ASE '15, 2015, pp. 295–306.
- [35] E. K. Smith, E. Barr, C. Le Goues, and Y. Brun, "Is the cure worse than the disease? Overfitting in automated program repair," in *Joint Meeting of the European Software Engineering Conference and the Symposium on the Foundations of Software Engineering*, 2015, pp. 532–543.
- [36] X. D. Le, F. Thung, D. Lo, and C. Le Goues, "Overfitting in semantics-based automated program repair," *Empirical Software Engineering*, vol. 23, no. 5, pp. 3007–3033, 2018.
- [37] C. S. Timperley, S. Stepney, and C. Le Goues, "An investigation into the use of mutation analysis for automated program repair," in *International Symposium on Search Based Software Engineering*, ser. SSBSE '17, 2017, pp. 99–114.
- [38] L. Beadle and C. Johnson, "Semantic analysis of program initialisation in genetic programming," *Genetic Programming and Evolvable Machines*, vol. 10, pp. 307–337, 2009.
- [39] Z. Lin, D. Marinov, H. Zhong, Y. Chen, and J. Zhao, "JaConTeBe: A benchmark suite of real-world Java concurrency bugs," in *Conference on Automated Software Engineering*, ser. ASE '15, 2015.
- [40] R. K. Saha, Y. Lyu, W. Lam, H. Yoshida, and M. R. Prasad, "BugsJar: A large-scale, diverse dataset of real-world Java bugs," in *International Conference on Mining Software Repositories*, ser. MSR '18, 2018, pp. 10–13.