

Test-equivalence Analysis for Automatic Patch Generation

SERGEY MECHTAEV, National University of Singapore

XIANG GAO, National University of Singapore

SHIN HWEI TAN, National University of Singapore

ABHIK ROYCHOUDHURY, National University of Singapore

Automated program repair is a problem of finding a transformation (called a patch) of a given incorrect program that eliminates the observable failures. It has important applications such as providing debugging aids, automatically grading student assignments and patching security vulnerabilities. A common challenge faced by existing repair techniques is scalability to large patch spaces, since there are many candidate patches that these techniques explicitly or implicitly consider.

The correctness criteria for program repair is often given as a suite of tests. Current repair techniques do not scale due to the large number of test executions performed by the underlying search algorithms. In this work, we address this problem by introducing a methodology of patch generation based on a test-equivalence relation (if two programs are "test-equivalent" for a given test, they produce indistinguishable results on this test). We propose two test-equivalence relations based on runtime values and dependencies respectively and present an algorithm that performs on-the-fly partitioning of patches into test-equivalence classes.

Our experiments on real-world programs reveal that the proposed methodology drastically reduces the number of test executions and therefore provides an order of magnitude efficiency improvement over existing repair techniques, without sacrificing patch quality.

CCS Concepts: • **Software and its engineering** → **Automatic programming; Software testing and debugging; Dynamic analysis;**

Additional Key Words and Phrases: Program repair; Program synthesis; Dynamic program analysis

ACM Reference Format:

Sergey Mechtaev, Xiang Gao, Shin Hwei Tan, and Abhik Roychoudhury. 2017. Test-equivalence Analysis for Automatic Patch Generation. *ACM Trans. Softw. Eng. Methodol.* 1, 1, Article 1 (January 2017), 37 pages. <https://doi.org/0000001.0000001>

1 INTRODUCTION

As every developer knows, debugging is difficult and extremely time-consuming. Due to the slow adoption of automated verification and debugging techniques, finding and eliminating defects remains mostly a manual process. Automated patch generation approaches can potentially alleviate this problem since they have been shown to be able to address defects in real-world programs and require minimal developer involvement. Specifically, they have been successfully applied for providing debugging hints [39], automatically grading assignments [30, 45] and patching security vulnerabilities [22]. However, the problem of huge search spaces pose serious challenges for current program repair techniques.

Authors' addresses: Sergey Mechtaev, National University of Singapore, mechtaev@comp.nus.edu.sg; Xiang Gao, National University of Singapore, gaoxiang@comp.nus.edu.sg; Shin Hwei Tan, National University of Singapore, shinhwei@comp.nus.edu.sg; Abhik Roychoudhury, National University of Singapore, abhik@comp.nus.edu.sg.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2017 Association for Computing Machinery.

1049-331X/2017/1-ART1 \$15.00

<https://doi.org/0000001.0000001>

The goal of program repair is to modify a given incorrect program to eliminate the observable failures. Specifically, the goal of *test-driven* program repair is to modify the buggy program so that it passes all given tests. The tests are typically provided by developers, but can also be automatically generated in a counterexample-guided refinement loop [1] when a formal specification is available. Patches (program modifications) that pass all given tests are referred to as *plausible* in the program repair literature [28]. Since a test-suite is an incomplete specification, plausible patches may not coincide with user’s intentions but may merely *overfit* the tests [32]. To address this problem, state-of-the-art techniques define a cost function (priority) on the space of candidate patches and search for a patch that optimizes this function. For example, changes can be prioritized based on syntactic distance [21], semantic distance [3] and information learned from human patches [17].

Patch generation systems need to consider large spaces of possible modifications in order to address many kinds of defects. One of the key challenges of program repair is scalability to large search spaces. Current techniques may require substantial time to generate patches and yet they consider and generate only relatively simple program transformations [15]. This impacts the ability of program repair to produce human-like repairs, since human patches often involve complex source code modifications. Besides that, a recent study [16] demonstrated that extending the search space with more transformations may cause repair systems to find fewer correct repairs because of the increased search time.

Although existing test-driven program repair techniques employ different methodologies (e.g. GenProg [42] uses genetic programming, SemFix [25] and Angelix [22] are based on constraint solving, Prophet [17] is based on machine learning), they all search for patches by repeatedly executing tests. Due to the high cost of test executions in large real-world programs, the number of performed test executions is the main bottleneck of many existing program repair algorithms.

Existing search methodologies can be divided into two categories: syntax-based and semantics-based. Syntax-based techniques (e.g. GenProg) explicitly generate and test syntactic changes. Thus, the number of test executions performed by such techniques is proportional to the number of explored candidate patches. Semantics-based techniques infer specification for identified statements through path exploration and synthesize changes based on this specification. For instance, Angelix explores deviations of execution paths for given tests using symbolic execution (which can be considered as a variant of test execution), and Prophet explores deviations of execution paths for given tests by enumerating sequences of condition values. Thus, the number of test executions performed by such techniques is proportional to the number of explored paths.

The purpose of this work is to improve the scalability of program repair without sacrificing the quality of generated patches. In order to achieve this, we propose a methodology based on a test-equivalence relation [8, 9]. If two programs are test-equivalent for a test, then the programs produce indistinguishable results on that test:

Definition 1.1 (Test-equivalence). Let \mathcal{P} be a set of programs, t be a test. An equivalence relation (reflexive, symmetric and transitive) $\sim^t \subset \mathcal{P} \times \mathcal{P}$ is a *test-equivalence relation* for t if it is consistent with the results of t , that is $\forall p_1, p_2 \in \mathcal{P}$, if $p_1 \sim^t p_2$ then p_1 and p_2 either both pass t or both fail t .

The proposed algorithm partitions the space of candidate patches into test-equivalence classes by performing on-the-fly analysis during test execution. This enables our methodology to alleviate the limitations of previous techniques. Compared with syntax-based techniques, it reduces the number of test executions since a single execution is sufficient to evaluate multiple patch candidates (specifically, all patches in the same test-equivalence class). Compared with semantics-based techniques, it reduces the number of test executions for two reasons. First, it avoids exploration of “infeasible” paths (sequences of values), i.e. paths or sequences of values that cannot be induced by any of the considered candidate patches in the context of given tests. Second, it reuses information inferred

across multiple tests to skip redundant executions, while previous semantics-based techniques perform path exploration independently for each test.

Contributions. The main contributions of this work are described in the following.

- (1) We propose the use of test-equivalence relations to drastically prune the search space explored for the purpose of program repair.
- (2) We adapt a test-equivalence relation used in mutation testing [8] (that we call value-based test-equivalence) to program repair by integrating it with syntax-guided program synthesis.
- (3) We define a new test-equivalence relation based on dynamic data dependencies and propose a method of synthesizing assignment statements through the composition of value-based and dependency-based test-equivalence relations.
- (4) We introduce a new patch space exploration algorithm that performs on-the-fly (during test execution) partitioning of patches into test-equivalence classes, thereby achieving efficient program repair that requires fewer test executions to generate a patch.
- (5) We conduct an evaluation of the algorithm on real-world programs from the GenProg ICSE'12 benchmark [12]; it demonstrates that test-equivalence significantly reduces the number of required test executions and therefore increases the efficiency of test-driven program repair and scales it to larger search spaces without sacrificing patch quality.

Outline. In the next section, we provide examples demonstrating limitations of existing techniques and formulate key insights of our methodology. Section 3 formally defines the two test-equivalence relations. Section 4 introduces a repair algorithm based on these relations, Section 5 describes its implementation and Section 6 presents its experimental evaluation. Section 7 discusses related work, Section 8 discusses future research directions and Section 9 concludes.

2 MOTIVATING EXAMPLES

This section gives three examples demonstrating limitations of existing techniques: a large number of redundant test executions, ineffectiveness in searching for optimal repairs and restricted applicability. We also formulate key insights that enable our method to address these limitations.

2.1 Example: repairing conditions

Consider a defect in the revision 0661f81 of Libtiff¹ from the GenProg ICSE'12 benchmark. The code in Figure 1a is responsible for flushing data written by the compression algorithm, and the defect is caused by the wrong highlighted condition. Libtiff test-suite contains 78 tests, and this defect is manifested by a failing test called “tiffcp-split”. Figure 1b demonstrates the developer patch that modifies the wrong condition by removing the clause `tif->tif_rawcc != orig_rawcc`.

We demonstrate how existing automated program repair algorithms generate a patch for this condition. First, repair algorithms perform fault localization to identify suspicious program statements. The number of localized statements in existing tools may vary from tens to thousands depending on algorithms and configurations (it can potentially include all executed statements). In this example, we consider only the location of the buggy expression highlighted in Figure 1a.

Second, program repair algorithms define a *search space* of candidate patches. In this work, we primarily focus on two state-of-the-art approaches that have been shown to scale to large real-world programs: Angelix [22] and Prophet [17]. Specifically, our goal was to support a combination of transformations implemented in these systems. Thus, the search space for the highlighted condition includes all possible replacements of its subexpressions by expressions constructed from visible program variables and C operators, refinements (e.g. appending `&& EXPR` and `|| EXPR`), replacements

¹Libtiff is a software library that provides support for TIFF image format: <http://simplesystems.org/libtiff/>

<pre> ... (*tif->tif_close)(tif); if (tif->tif_rawcc > 0 && tif->tif_rawcc != orig_rawcc && (tif->tif_flags & TIFF_BEENWRITING) != 0 && !TIFFFlushData1(tif)) { TIFFErrorExt(tif->tif_clientdata, module, "Error_flushing_data_before_ directory_write"); return (0); } ... </pre> <p>(a) Incorrect condition in Libtiff (rev. 0661f81).</p>	<pre> ... (*tif->tif_close)(tif); if (tif->tif_rawcc > 0 && (tif->tif_flags & TIFF_BEENWRITING) != 0 && !TIFFFlushData1(tif)) { TIFFErrorExt(tif->tif_clientdata, module, "Error_flushing_data_before_ directory_write"); return (0); } ... </pre> <p>(b) Developer patch for incorrect condition.</p>
--	---

Fig. 1. Defect in Libtiff library from GenProg ICSE'12 benchmark.

of operators and swapping arguments. In total, the search space in our synthesizer contains 56 243 modifications of the buggy condition.

Finally, program repair algorithms explore the search space in order to try to find a modification that passes all given tests. We say that an element of a search space is *explored* if the algorithm identifies if it passes all the tests or fails at least one. Existing search space exploration methods can be classified into two categories: syntax-based and semantics-based. *Syntax-based* algorithms explicitly generate and test syntactic changes. In this example, a syntax-based algorithm have to execute the failing test 56 243 times to evaluate all candidates². Since there are 78 tests in the test-suite, 907 457 test executions are required to explore the search space³. Given the high cost of test execution, this approach has poor scalability.

Semantics-based techniques (e.g. Semfix [25], SPR [15], Angelix and Prophet) split exploration into two phases. First, they infer a synthesis specification for the identified expression through path exploration. For this example, they enumerate and execute sequences of condition values (e.g. *true, true, true, false, ...*) to find those sequences that enable the program to pass the test. Second, they synthesize a modification of the condition to match the inferred specification. In this example, there are 256 possible execution paths (the condition is evaluated multiple times during the test execution), therefore a semantics-based algorithm performs 256 test executions for the failing tests, and 1320 for the whole test-suite⁴. Although semantics-based techniques were shown to be more scalable [15], they are subject to the *path explosion problem*: the number of execution paths can be infinite. To address this, current systems introduce a bound for the number of explored paths, however it may affect their effectiveness: if a path followed by the correct patch is omitted, then this correct patch cannot be generated.

The algorithm proposed in this work performs on-the-fly partitioning of program modifications into test-equivalence classes. We demonstrate the effect of the relation $\overset{t}{\sim}_{value}$ described in Section 3.3. Two modifications of a program expression are test-equivalent w.r.t. $\overset{t}{\sim}_{value}$ if they are evaluated into the same sequences of values during the test execution. Surprisingly, the space of 56 243

²Since the search space contains the correct patch in this example, the algorithm can stop search earlier after the patch is found. Then, the number of test executions depends on the exploration order.

³This data is obtained by executing our implementation of syntactic enumeration.

⁴This data is obtained by executing Angelix.

1. `((tif->tif_rawcc > 0) && (tif->tif_rawcc != orig_rawcc))
 || (tif->tif_flags & TIFF_BEENWRITING)`
2. `((tif->tif_rawcc > 0) || (tif->tif_rawcc != orig_rawcc))
 && (tif->tif_flags & TIFF_BEENWRITING)`
3. `((tif->tif_rawcc == 0) && (tif->tif_rawcc != orig_rawcc))
 && (tif->tif_flags & TIFF_BEENWRITING)`
4. `((tif->tif_rawcc > 0) && (tif->tif_rawcc != orig_rawcc))
 && (tif->tif_flags & TIFF_BEENWRITING) || (imagedone >= orig_rawcc)`
5. `((tif->tif_rawcc > 0) && (tif->tif_rawcc != orig_rawcc))
 && (tif->tif_flags & TIFF_BEENWRITING) || (tif->tif_flags >= 74)`

Fig. 2. Each of 56 243 search space elements is test-equivalent to one of these 5 expressions.

modifications can be partitioning into only 5 test-equivalence classes for the failing test “tiffcp-split” w.r.t. $\overset{t}{\sim}_{value}$; five elements of the search space that represent different test-equivalence classes are given in Figure 2. Since all patches in the same test-equivalence class exhibit the same behaviour for the corresponding test, the failing test can be executed only 5 times to evaluate all candidates.

Our algorithm computes test-equivalence classes for each test in the test-suite. However, since test-equivalence classes for different tests may intersect, our algorithm takes advantage of this to skip redundant execution across different tests. Specifically, for each next test it only evaluates subspaces of modifications that are not included into failing test-equivalence classes of previously executed tests. Meanwhile, semantics-based techniques perform specification inference for each test independently without reusing information across tests. As a result, our algorithm requires only 103 test executions to evaluate all 56 243 modifications with the whole test-suite.

Key insight. The key insight that enables our method to reduce the number of required test executions is that, compared with techniques that explore execution paths, it takes the expressiveness of the patch space into account (e.g. it identifies that only 5 out of 256 possible execution paths are induced by the considered set of 56 243 transformations). Compared with syntactic enumeration, it substantially reduces executions since a single execution evaluates a whole test-equivalence class.

2.2 Example: optimal repair

Since a test-suite is an incomplete specification, test-driven program repair suffers from the test overfitting problem [32]. To address this issue, state-of-the-art techniques define a priority (a cost function) in the space of patches and search for a program modification that optimizes this function. Ideally, this function should assign higher cost to overfitting patches. For instance, Prophet [17] demonstrates how such a cost function learned from human patches enables the generation of more correct repairs.

Consider a program p in Figure 3a that counts odd numbers in the interval $(0, i]$. The $*$ indicates a wrong condition that has to be modified by the repair algorithm (the correct condition is $i \bmod 2 = 1$). We denote a program obtained by substituting $*$ with an expression e as $p[*/e]$. The repair algorithm searches for a plausible patch (a substitution of $*$ with a condition) from the space \mathcal{P} in Figure 3b such that the resulting program passes the test t defined as follows:

$$t := (\{ i \mapsto 4, c \mapsto 0 \}, \lambda\sigma. \sigma(c) = 2)$$

where t is pair of (1) an initial program state (mapping from variables to values) and (2) a test assertion (a boolean function over program states) denoted using lambda notation. We assume that $*$ is such that p fails t . Besides that, we consider a cost function κ defined for the considered space of substitutions in Figure 3c. The goal is to find a plausible patch with the lowest cost.

<pre> while i > 0 do if * then c := c + 1 fi; i := i - 1 od </pre> <p>(a) Buggy program p.</p>	$\mathcal{P} := \{ p[*i \geq 0],$ $p[*c \geq 0],$ $p[*i \bmod 2 = 1],$ $p[*i \bmod 2 = 0],$ $p[*i > 2] \}$ <p>(b) Search space.</p>	$\kappa(p[*i \geq 0]) := 0.1$ $\kappa(p[*c \geq 0]) := 0.2$ $\kappa(p[*i \bmod 2 = 1]) := 0.3$ $\kappa(p[*i \bmod 2 = 0]) := 0.4$ $\kappa(p[*i > 2]) := 0.5$ <p>(c) Cost function.</p>
--	---	--

Fig. 3. Example of optimal program repair problem.

In order to find a patch for the example program, techniques like Angelix and Prophet enumerate possible sequences of values that a condition can take during test execution. Since there can be potentially infinite number of such sequences, existing approaches introduce a bound for the number of explored sequences and use an exploration heuristics to choose which sequences to explore. For instance, Prophet enumerates sequences where the condition first always takes the true branch until a certain point after which it always takes the false branch. Thus, for the considered example it would enumerate the following sequences:

$$\{ true, true, true, true \},$$

$$\{ true, true, true, false \},$$

$$\{ true, true, false, false \},$$

$$\{ true, false, false, false \}$$

For each of these sequences, Prophet executes the program with the test t in such a way that the condition $*$ takes the values as in this sequence during the execution. Only the third sequence $\{ true, true, false, false \}$ enables the program to pass t , therefore it will be selected as a specification for expression synthesis. The synthesizer will find the expression $i > 2$ obtaining a suboptimal patch $p[*i > 2]$ with the cost 0.5, since this is the only expression from the search space satisfying the specification. However, the correct expression $i \bmod 2 = 1$ with a lower cost 0.3 cannot be generated, since the corresponding sequence $\{ false, true, false, true \}$ is not explored by the algorithm.

In contrast to techniques like Angelix and Prophet, our algorithm iterates through the search space in such a way that at each steps it selects and evaluates an unevaluated candidate with the lowest cost. Specifically, it starts by choosing the candidate $p[*i \geq 0]$ with the cost 0.1. It executes this candidate on-the-fly computing its test-equivalence class w.r.t. $\overset{t}{\sim}_{value}$ described in Section 2.1. This class contains the program $p[*c \geq 0]$, since the conditions $i \geq 0$ and $c \geq 0$ produce the same sequence of values $\{ true, true, true, true \}$ for t . Since $p[*i \geq 0]$ does not pass the test, the whole corresponding test-equivalence class is marked as failing. Next, it selects $p[*i \bmod 2 = 1]$ with the cost 0.3 since $p[*c \geq 0]$ was indirectly evaluated through test-equivalence at the previous step. Since this candidate passes the test, the algorithm outputs it as a found repair.

Key insight. Our algorithm guides exploration based on a given cost function and focuses on high priority areas of the space of patches. By construction, if it finds a patch, then this patch is guaranteed to be the global optimum in the search space w.r.t. the cost function. Angelix and Prophet, on the other hand, may spend executions for value sequences that correspond to suboptimal candidates or correspond to no candidates at all (e.g. $\{ false, true, true, true \}$), and therefore may miss the best patch in their search space.

<pre> ... clear_bufs(); to_stdout = 1; part_nb = 0; ifd = part_nb; if (decompress) { method=get_method(ifd); ... </pre>	<pre> ... clear_bufs(); to_stdout = 1; part_nb = 0; ifd = 0; if (decompress) { method=get_method(ifd); ... </pre>	<pre> ... clear_bufs(); to_stdout = 1; part_nb = 0; if (decompress) { ifd = part_nb; method=get_method(ifd); ... </pre>
(a) Before if-statement.	(b) Before if-statement.	(c) Inside if-statement.

Fig. 4. Candidate patches for defect of Gzip from GenProg ICSE'12 benchmark.

2.3 Example: repairing assignments

Although current program repair approaches have been shown to be relatively effective in modifying existing program expressions, they provide limited support for more complex transformations. In this work, we consider one such transformation that inserts a new assignment statement to the buggy program. Techniques like Prophet and GenProg can generate patches by copying/moving existing program assignments, however this approach has limitations: (1) assignments for local variables cannot be copied from different parts of the program because of their scope and (2) each insertion of an assignment is validated separately, which yield a large number of required test executions. Existing techniques do not apply specification inference for assignment synthesis (as described in Section 2.1 for conditions) because such specification has to encode all possible side effects that can be caused by assignment insertion (for each variable that can appear in the left-hand side of the assignment), which makes inferring such specification infeasible for large programs.

We show how test-equivalence can scale assignment synthesis for a defect in Gzip⁵ from the GenProg ICSE'12 benchmark. Consider three candidate patches in Figure 4 that insert the highlighted statements at several program locations. First, our algorithm identifies that the program in Figure 4a is test-equivalent to the program in Figure 4b (w.r.t. the relation \sim_{value}^t described previously in Section 2.1) since they differ only in the right-hand side of the highlighted assignments and the corresponding expressions take the same values during test execution. Second, using a simple dynamic data dependency analysis our algorithm identifies that the program in Figure 4a is test-equivalent to the program in Figure 4c since (1) they insert the same assignment at different program locations, (2) both these locations are executed by the test since the true branch of the if-statement is taken during the test execution and (3) the variables `ifd` and `part_nb` are not used/modified between these locations during test execution. We refer to such a test-equivalence relation as \sim_{deps}^t . Finally, our algorithm merges the results of the two analyses (as the transitive closure of their union) and determines that the program in Figure 4b is test-equivalent to the program in Figure 4c. Therefore, a single test execution is sufficient to evaluate all these patches.

Key insight. Since test-equivalence is a weaker property than the property of “passing the test” expressed by the inferred specification in semantics-based techniques, it permits using more lightweight analysis techniques. Specifically, we demonstrate that a composition of two lightweight test-equivalence analyses enables us to scale assignment synthesis.

⁵Gzip is a file compression/decompression application: <https://www.gnu.org/software/gzip/>

$\langle Stmt \rangle ::= \langle Var \rangle := \langle AExpr \rangle$ $\quad \text{ skip}$ $\quad \langle Stmt \rangle ; \langle Stmt \rangle$ $\quad \text{ if } \langle BExpr \rangle \text{ then } \langle Stmt \rangle \text{ else } \langle Stmt \rangle \text{ fi}$ $\quad \text{ while } \langle BExpr \rangle \text{ do } \langle Stmt \rangle \text{ od}$ $\langle AExpr \rangle ::= \langle Var \rangle$ $\quad \langle Num \rangle$ $\quad \langle AExpr \rangle \langle AOp \rangle \langle AExpr \rangle$	$\langle BExpr \rangle ::= \text{ true}$ $\quad \text{ false}$ $\quad \langle BExpr \rangle \langle BOp \rangle \langle BExpr \rangle$ $\quad \langle AExpr \rangle \langle ROp \rangle \langle AExpr \rangle$ $\langle AOp \rangle ::= + - * \dots$ $\langle BOp \rangle ::= \text{ and } \text{ or } \dots$ $\langle ROp \rangle ::= < <= = \dots$
--	--

Fig. 5. Syntax of programming language \mathcal{L} .

3 TEST-EQUIVALENCE RELATIONS

This section formally introduces two test-equivalence relations for spaces of program modifications generated through program synthesis. In the subsequent Section 4, we demonstrate how these relations can be applied for scaling patch generation. However, we believe that these relations can be also used in different domains; other potential applications are discussed in Section 7.

3.1 Preliminaries

We introduce our methodology for an imperative programming language \mathcal{L} . The syntax of \mathcal{L} is defined in Figure 5, where \mathbb{Z} is the integer domain, \mathbb{B} is the boolean domain (true and false), $Stmt$ is a set of statements, $AExpr$ is a set of arithmetic expressions, $BExpr$ is a set of boolean expressions, $Expr = AExpr \cup BExpr$, Num is a set of integer literals, Var is a set of variables over \mathbb{Z} . A program in \mathcal{L} is a sequence of statements. We denote subsets of \mathcal{L} as \mathcal{P} , subsets of expressions $Expr$ as \mathcal{E} , all variables from Var encountered in an expression e as $Var(e)$, a program that is obtained by substituting a statement (expression) s with a statement (expression) s' in a program p as $p[s/s']$.

The semantics of \mathcal{L} is defined in Figure 6, where program state $\sigma : Var \rightarrow \mathbb{Z}$ is a function from program variables into values, Σ is a set of program states. We indicate a modification of a program state σ where the value of the variable v is updated to n as $\sigma[v \mapsto n]$.

Definition 3.1 (Test). Let $p \in \mathcal{L}$ be a program, $t \in \Sigma \times (\Sigma \rightarrow \mathbb{B})$ be a test, that is a pair (σ_{in}, ϕ) , where σ_{in} is the initial program state (input) and ϕ is the test assertion (a boolean function over program states). We say that p passes t (indicated as $Pass[p, t]$) iff $\langle p, \sigma_{in} \rangle \Downarrow \sigma_{out} \wedge \phi(\sigma_{out})$.

A test execution can be represented through a derivation tree as in the following example:

Example 3.2 (Derivation tree). Consider a program p defined as “ $x := y + 1; y := x$ ” and a test $t := (\sigma_{in}, \lambda\sigma. \sigma(y) = 3)$ where the input state $\sigma_{in} := \{x \mapsto 1, y \mapsto 2\}$. According to the semantics in Figure 6, the following relation holds: $\langle p, \sigma_{in} \rangle \Downarrow \sigma_{out}$, where $\sigma_{out} := \{x \mapsto 3, y \mapsto 3\}$. This relation can be established by the following *derivation tree* obtained by applying the semantics rules:

$$\frac{\frac{\frac{\overline{\langle y, \dots \rangle \Downarrow 2} \text{VAR}}{\langle y + 1, \{x \mapsto 1, y \mapsto 2\} \Downarrow 3} \text{OP}}{\langle x := y + 1, \{x \mapsto 1, y \mapsto 2\} \Downarrow \{x \mapsto 3, y \mapsto 2\} \text{ASSIGN}} \quad \frac{\overline{\langle 1, \dots \rangle \Downarrow 1} \text{NUM}}{\langle x, \dots \rangle \Downarrow 3 \text{VAR}}}{\langle y := x, \{x \mapsto 3, y \mapsto 2\} \Downarrow \{x \mapsto 3, y \mapsto 3\} \text{ASSIGN}}}{\langle x := y + 1; y := x, \{x \mapsto 1, y \mapsto 2\} \Downarrow \{x \mapsto 3, y \mapsto 3\} \text{SEQ}}$$

The output state $\{x \mapsto 3, y \mapsto 3\}$ satisfies the test assertion $\lambda\sigma. \sigma(y) = 3$, therefore p passes t .

VAR	NUM	OP	ASSIGN
$\langle v, \sigma \rangle \Downarrow \sigma[v]$	$\langle n, \sigma \rangle \Downarrow n$	$\langle e_1, \sigma \rangle \Downarrow n_1 \quad \langle e_2, \sigma \rangle \Downarrow n_2 \quad n_3 = n_1 \text{ op } n_2$ $\langle e_1 \text{ op } e_2, \sigma \rangle \Downarrow n_3$	$\langle e, \sigma \rangle \Downarrow n$ $\langle v := e, \sigma \rangle \Downarrow \sigma[v \mapsto n]$
SEQ	IF-TRUE	IF-FALSE	
$\langle s_1, \sigma \rangle \Downarrow \sigma_1 \quad \langle s_2, \sigma_1 \rangle \Downarrow \sigma_2$ $\langle s_1 ; s_2, \sigma \rangle \Downarrow \sigma_2$	$\langle e, \sigma \rangle \Downarrow \text{true} \quad \langle s_1, \sigma \rangle \Downarrow \sigma_1$ $\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \Downarrow \sigma_1$	$\langle e, \sigma \rangle \Downarrow \text{false} \quad \langle s_2, \sigma \rangle \Downarrow \sigma_2$ $\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \Downarrow \sigma_2$	
WHILE-TRUE	WHILE-FALSE	SKIP	
$\langle e, \sigma \rangle \Downarrow \text{true} \quad \langle s_1, \sigma \rangle \Downarrow \sigma_1$ $\langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \Downarrow \sigma_2$	$\langle e, \sigma \rangle \Downarrow \text{false}$ $\langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \Downarrow \sigma$	$\langle \text{skip}, \sigma \rangle \Downarrow \sigma$	

Fig. 6. Semantics of \mathcal{L} . v – variables, n – integer values, e – expressions, s – statements, σ – program states.**ALGORITHM 1:** Value-projection operator via enumerative synthesis**Input:** set of expressions \mathcal{E} , program state σ , value n **Output:** set of expressions \mathcal{E}'

```

1  $\mathcal{E}' := \emptyset;$ 
2 foreach  $e \in \mathcal{E}$  do
3   if  $\langle e, \sigma \rangle \Downarrow n$  then
4      $\mathcal{E}' := \mathcal{E}' \cup \{e\};$ 
5 return  $\mathcal{E}';$ 

```

3.2 Generalized synthesis

The presented methodology of test-equivalence analysis is introduced for program repair techniques that rely on syntax-guided program synthesis to generate patches. Such techniques generate expressions for “holes” in programs based on the values of visible program variables (program state) and expected results for the synthesized expressions. Thus, we define *synthesis specification* as a finite set of input-output pairs (an input is represented by a program state from Σ and an output is an integer or a boolean value); we denote the set of all specifications as $\text{Spec} := 2^{\Sigma \times (\mathbb{Z} \cup \mathbb{B})}$.

Definition 3.3 (Synthesis procedure). A syntax-guided synthesis procedure $\text{synthesize} : 2^{\text{Expr}} \times \text{Spec} \rightarrow \text{Expr}$ is a function that takes a set of expressions (the synthesis search space) and a specification, and returns an expression from the search space that meets the specification. Specifically, for a given search space \mathcal{E} and specification spec , if $\text{synthesize}(\mathcal{E}, \text{spec}) = e$ then $e \in \mathcal{E} \wedge \bigwedge_{\sigma, n \in \text{spec}} \langle e, \sigma \rangle \Downarrow n$.

In order to integrate synthesis with test-equivalence analysis, we impose additional requirements for the program synthesizer: it should define a value-projection operator over its search space.

The *value-projection operator* $\Pi_{\sigma, n}^{\text{value}}$ produces a maximal subset of a given set of expressions consisting only of expressions that are evaluated into n in the context σ :

$$\Pi_{\sigma, n}^{\text{value}}(\mathcal{E}) = \{e \mid e \in \mathcal{E} \wedge \langle e, \sigma \rangle \Downarrow n\}$$

In this work, we use an enumerative synthesizer [1] that demonstrated positive results in program synthesis competitions⁶. Since it represents the search space explicitly as a set of expressions, it is straightforward to realize the value-projection operator in such a synthesizer (Algorithm 1). Other possible realizations are discussed in Section A.

<p>EXPR-MOD</p> $\frac{\text{Modified}(e) \quad \langle e, \sigma \rangle \Downarrow n \quad c' = \Pi_{\sigma, n}^{\text{value}}(c)}{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle n, c' \rangle}$	<p>EXPR-NMOD</p> $\frac{\neg \text{Modified}(e) \quad \langle e, \sigma \rangle \Downarrow n}{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle n, c \rangle}$	<p>ASSIGN</p> $\frac{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle n, c' \rangle}{\langle v := e, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma[v \mapsto n], c' \rangle}$
<p>IF-TRUE</p> $\frac{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle \text{true}, c_1 \rangle \quad \langle s_1, \sigma, c_1 \rangle \Downarrow_{\text{value}} \langle \sigma_1, c_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma_1, c_2 \rangle}$	<p>IF-FALSE</p> $\frac{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle \text{false}, c_1 \rangle \quad \langle s_2, \sigma, c_1 \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_2 \rangle}$	
<p>SEQ</p> $\frac{\langle s_1, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma_1, c_1 \rangle \quad \langle s_2, \sigma_1, c_1 \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_2 \rangle}{\langle s_1 ; s_2, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_2 \rangle}$	<p>SKIP</p> $\frac{}{\langle \text{skip}, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma, c \rangle}$	<p>WHILE-FALSE</p> $\frac{}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma, c' \rangle}$
<p>WHILE-TRUE</p> $\frac{\langle e, \sigma, c \rangle \Downarrow_{\text{value}} \langle \text{true}, c_1 \rangle \quad \langle s_1, \sigma, c_1 \rangle \Downarrow_{\text{value}} \langle \sigma_1, c_2 \rangle \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma_1, c_2 \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_3 \rangle}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, c \rangle \Downarrow_{\text{value}} \langle \sigma_2, c_3 \rangle}$		

Fig. 7. Augmented semantics of \mathcal{L} for computing test-equivalence classes w.r.t. $\overset{t}{\sim}_{\text{value}}$. v – variables, n – integer values, e – expressions, s – statements, σ – program states, c – sets of expressions, *Modified* – predicate over expressions.

3.3 Value-based test-equivalence relation

This section introduces a test-equivalence relation $\overset{t}{\sim}_{\text{value}}$ for spaces of programs that differ only in expressions. From the conceptual point of view, the introduced relation generalizes relations that have been used in mutation testing [8]. From the algorithmic point of view, our approach differs from the previous work on test-equivalence in that the relation is integrated with syntax-guided program synthesis (used for patch synthesis) via the value-projector operator.

Intuitively, two programs p and p' such that $p' = p[e/e']$ for some expressions e and e' are test-equivalent for some test t w.r.t. $\overset{t}{\sim}_{\text{value}}$ if, during the executions of p and p' with t , the expressions e and e' are evaluated into the same values. An example of applying this relation is given in Section 2.1.

We define the relation $\overset{t}{\sim}_{\text{value}}$ *constructively* using an augmented semantics of \mathcal{L} . We chose this presentation since it simultaneously defines an algorithm of computing test-equivalence classes in spaces of modification generation through program synthesis. The implementation of this semantics via program instrumentation is discussed in Section 5.

The semantics in Figure 7 extends the semantics in Figure 6 by defining the function $\Downarrow_{\text{value}}$. It is parameterized by a predicate $\text{Modified} : \text{Expr} \rightarrow \mathbb{B}$ that marks the modified program expression, substitutions of which are analyzed for test-equivalence. The function $\Downarrow_{\text{value}}$ additionally maintains a set of expressions (denoted as c), such that the substitutions of the modified expression with c form the computed test-equivalent class.

The augmented semantics describes an algorithm of identifying test-equivalence classes that, for a given set of expressions (the synthesis search space), “filters out” those that do not belong to the test-equivalent class of the current program by repeatedly applying the value-projection operator. The application of the value-projection operator to the current set of expressions c is highlighted in Figure 7. Thus, it identifies all expressions from c that produce the same value as the original expression at this evaluation step. Since each expression can be evaluated multiple times during test execution, the value-projection operator can also be applied multiple times. Therefore, the test-equivalence class is computed as $\Pi_{\sigma_1, n_1}^{\text{value}} \circ \Pi_{\sigma_2, n_2}^{\text{value}} \circ \dots \circ \Pi_{\sigma_k, n_k}^{\text{value}}(\mathcal{E})$, where \mathcal{E} is a set of all

⁶SyGuS-Comp 2014: <http://www.syguus.org/SyGuS-COMP2014.html>

substitutions of the modified expression, n_i are the values of the modified expression computed during test execution and σ_i are the corresponding program states.

Definition 3.4 (Value-based test-equivalence relation). Let \mathcal{P} be a set of programs, $t = (\sigma_{in}, \phi)$ be a test. $\overset{t}{\sim}_{value} \subset \mathcal{P} \times \mathcal{P}$ is a *value-based test-equivalence relation* iff $p_1 \overset{t}{\sim}_{value} p_2$ for $p_1, p_2 \in \mathcal{P}$ if $\exists e, e' \in Expr$ such that $p_2 = p_1[e/e']$ and $\langle p_1, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle _, \{e, e'\} \rangle$ given that $Modified := \lambda x. x = e$.

In this definition, we call two programs that differ only in expressions to be test-equivalent if the corresponding expressions produce the same values according to the semantics in Figure 6. Specifically, by passing the program p_1 , the test input σ_{in} and the set of expressions $\{e, e'\}$ as the arguments to \Downarrow_{value} , we obtain the same set $\{e, e'\}$ as the result.

PROPOSITION 3.5. *The relation $\overset{t}{\sim}_{value}$ is a test-equivalence relation according to Definition 1.1.*

The proposition above formally states that (1) $\overset{t}{\sim}_{value}$ is an equivalence relations and (2) if two programs that differ only in expressions are test-equivalent according to the semantics in Figure 7, then these two programs either both pass the test or both fail the test. A proof for the above proposition is given in Section C.

Example 3.6 (Test-equivalent programs w.r.t. $\overset{t}{\sim}_{value}$). Consider a program p_1 defined as

if $x > 0$ then $x := y$ else skip fi

a program p_2 defined as

if $y = 2$ then $x := y$ else skip fi

and a test $t := (\sigma_{in}, \lambda \sigma. \sigma(y) = 3)$ where the input state $\sigma_{in} := \{x \mapsto 1, y \mapsto 2\}$. These programs are test-equivalent w.r.t. the value-based test-equivalence relation for the test t , since they differ only in the if-condition and the following relation holds:

$$\langle p_1, \sigma_{in}, \{“x > 0”, “y = 2”\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{“x > 0”, “y = 2”\} \rangle$$

where $\sigma_{out} := \{x \mapsto 2, y \mapsto 2\}$. This relation can be established by the following derivation tree:

$$\frac{\frac{\dots}{\langle x > 1, \sigma_{in} \rangle \Downarrow true} \text{OP} \quad \{“x > 0”, “y = 2”\} = \Pi_{\sigma_{in}, true}^{value} (\{“x > 0”, “y = 2”\})}{\langle x > 0, \sigma_{in}, \{“x > 0”, “y = 2”\} \rangle \Downarrow_{value} \langle true, \{“x > 0”, “y = 2”\} \rangle} \text{EXPR-MOD} \quad \frac{\dots}{\dots} \text{ASSIGN}}{\langle \text{if } x > 0 \text{ then } x := y \text{ else skip fi}, \sigma_{in}, \{“x > 0”, “y = 2”\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{“x > 0”, “y = 2”\} \rangle} \text{IF-TRUE}$$

3.4 Dependency-based test-equivalence relation

This section introduces a test-equivalence relation $\overset{t}{\sim}_{deps}$ for spaces of programs that differ in locations in which an assignment statement is inserted. Let a *location* in program p be a statement of p . We say that a program p' is obtained by inserting the assignment $v := e$ at the location l iff $p' = p[l/v := e, l]$. Let p be a program and programs p_1 and p_2 are obtained by inserting the assignment $v := e$ at the locations l_1 and l_2 of p respectively. Informally, p_1 and p_2 are test-equivalent for some test t if, during an execution of p_1 with t , (1) for each occurrence of l_1 in the execution trace there is a “matching” occurrence of l_2 (the variable v is not read or overwritten between these occurrences and the variables $Var(e)$ are not overwritten between these occurrences), and (2) for each occurrence of l_2 in the execution trace there is a “matching” occurrence of l_1 . An example of applying this relation is given in Section 2.3.

$\frac{\text{VAR-LEFT-EXE} \quad \text{Left}(v) \wedge x}{\langle v, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma(v), l \cap c, \emptyset, \text{false} \rangle}$ $\frac{\text{OP} \quad \langle e_1, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle n_1, l_1, c_1, x_1 \rangle}{\langle e_1 \text{ op } e_2, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle n_3, l_2, c_2, x_2 \rangle}$	$\frac{\text{VAR-LEFT-NEXE} \quad \text{Left}(v) \wedge \neg x}{\langle v, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma(v), l \setminus c, \emptyset, x \rangle}$ $\frac{\text{ASSIGN-INS} \quad \text{Inserted}(v := e) \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma', l, c, \text{true} \rangle}$ $\frac{\text{ASSIGN-LR-EXE} \quad \neg \text{Inserted}(v := e) \wedge (\text{Left}(v) \vee \text{Right}(v)) \wedge x \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma', l \cap c, \emptyset, \text{false} \rangle}$ $\frac{\text{ASSIGN-LR-NEXE} \quad \neg \text{Inserted}(v := e) \wedge (\text{Left}(v) \vee \text{Right}(v)) \wedge \neg x \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma', l \setminus c, \emptyset, \text{false} \rangle}$	$\frac{\text{VAR-NLEFT} \quad \neg \text{Left}(v)}{\langle v, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma(v), l, c, x \rangle}$ $\frac{\text{ASSIGN-NLR} \quad \neg \text{Inserted}(v := e) \wedge \neg \text{Left}(v) \wedge \neg \text{Right}(v) \quad \langle e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle n, l_1, c_1, x_1 \rangle}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma[v \mapsto n], l_1, c_1, x_1 \rangle}$
$\frac{\text{SEQ} \quad \langle s_1, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma_1, l_1, c_1, x_1 \rangle \quad \langle s_2, \sigma_1, l_1, c_1 \cup \{s_2\}, x_1 \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_2, c_2, x_2 \rangle}{\langle s_1 ; s_2, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_2, c_2, x_2 \rangle}$	$\frac{\text{IF-TRUE} \quad \langle e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \text{true}, l_1, c_1, x_1 \rangle \quad \langle s_1, \sigma, l_1, c_1 \cup \{s_1\}, x_1 \rangle \Downarrow_{\text{deps}} \langle \sigma_1, l_2, c_2, x_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, c \rangle \Downarrow_{\text{deps}} \langle \sigma_1, l_2, c_2, x_2 \rangle}$	$\frac{\text{SKIP} \quad \langle \text{skip}, \sigma, l, c, x \rangle}{\langle \text{skip}, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma, l, c, x \rangle}$
$\frac{\text{IF-FALSE} \quad \langle e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \text{false}, l_1, c_1, x_1 \rangle \quad \langle s_2, \sigma, l_1, c_1 \cup \{s_2\}, x_1 \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_2, c_2, x_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_2, c_2, x_2 \rangle}$	$\frac{\text{WHILE-TRUE} \quad \langle e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \text{true}, l_1, c_1, x_1 \rangle \quad \langle s_1, \sigma, l_1, c_1 \cup \{s_1\}, x_1 \rangle \Downarrow_{\text{deps}} \langle \sigma_1, l_2, c_2, x_2 \rangle \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma_1, l_2, c_2, x_2 \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_3, c_3, x_3 \rangle}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma_2, l_3, c_3, x_3 \rangle}$	$\frac{\text{WHILE-FALSE} \quad \langle e, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \text{false}, l', c', x' \rangle \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma, l', c', x' \rangle}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle \sigma, l', c', x' \rangle}$
$\frac{\text{NUM} \quad \langle n, \sigma \rangle \Downarrow n}{\langle n, \sigma, l, c, x \rangle \Downarrow_{\text{deps}} \langle n, l, c, x \rangle}$		

Fig. 8. Augmented semantics of \mathcal{L} for computing test-equivalence classes w.r.t. \downarrow_{deps} . v – variables, n – integer values, e – expressions, s – statements, σ – program states, l, c – sets of locations, x – boolean values, Inserted – predicate over statements, $\text{Left}, \text{Right}$ – predicates over variables.

The relation is formally defined through an augmented semantics of \mathcal{L} . As in Section 3.3, we chose this representation since it simultaneously defines an algorithm of identifying test-equivalence classes. The implementation of this semantics via program instrumentation is discussed in Section 5.

The semantics in Figure 8 extends the semantics in Figure 6 by defining the function \downarrow_{deps} . It is parameterized by a predicate $\text{Inserted} : \text{Stmt} \rightarrow \mathbb{B}$ that marks the inserted assignment, a predicate $\text{Left} : \mathcal{V} \rightarrow \mathbb{B}$ that marks the left-hand side variable of the inserted assignment, and a predicate $\text{Right} : \mathcal{V} \rightarrow \mathbb{B}$ that marks the variables used in the right-hand side of the inserted assignment. \downarrow_{deps} additionally maintains:

- l – a set of locations representing test-equivalent insertions;
- c – a set of locations that are executed after the last read/write of the variables involved in the inserted assignment;
- x – a boolean value that indicates if the inserted assignment was evaluated after the last read/write of the variables involved in the inserted assignment.

The augmented semantics describes an analysis algorithm that, for a given set of locations, “filters out” those that do not correspond to the test-equivalent insertions of a given assignment. For a

program with inserted assignment $v := e$, the semantics in Figure 8 computes sequences of executed locations (stored in the set c) such that each sequence contains the inserted statement $v := e$ (the rule ASSIGN-INS) and all the rest of the statements in this sequence do not read/overwrite the variable v and do not overwrite the variables in e (the rule ASSIGN-NLR). When such a sequence is found, the set of locations executed in this sequence is intersected with the current set of test-equivalent insertions ($l \cap c$ in the rules VAR-LEFT-EXE and ASSIGN-LR-EXE). When the inserted assignment $v := e$ is not executed in such a sequence, the set of locations executed in this sequence is removed from the set of test-equivalent insertions ($l \setminus c$ in the rules VAR-LEFT-NEXE and ASSIGN-LR-NEXE), since for these locations this is no “matching” occurrence of $v := e$.

Definition 3.7 (Dependency-based test-equivalence relation). Let \mathcal{P} be a set of programs, $t = (\sigma_{in}, \phi)$ be a test. $\overset{t}{\sim}_{deps} \subset \mathcal{P} \times \mathcal{P}$ is a *dependency-based test-equivalence relation* iff $p_1 \overset{t}{\sim}_{deps} p_2$ for $p_1, p_2 \in \mathcal{P}$ if there is program p with locations l_1, l_2 such that $p_1 = p[l_1/v := e; l_1]$ and $p_2 = p[l_2/v := e; l_2]$ and $\langle p_1, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{l_1, l_2\}, _, _ \rangle$ given that *Inserted* := $(\lambda s. s = “v := e”)$, *Left* := $(\lambda v'. v' = v)$, *Right* := $(\lambda v'. v' \in Var(e))$.

In this definition, we call two programs that differ in locations of an assignment insertion to be test-equivalent if the difference does not affect dynamic data dependencies according to the semantics in Figure 6. Specifically, by passing the program p_1 , the test input σ_{in} and the set of locations $\{l_1, l_2\}$ as the arguments to \Downarrow_{deps} , we obtain the same set $\{l_1, l_2\}$ as the result.

PROPOSITION 3.8. *The relation $\overset{t}{\sim}_{deps}$ is a test-equivalence relation according to Definition 1.1.*

The proposition above formally states that (1) $\overset{t}{\sim}_{deps}$ is an equivalence relations and (2) if two programs that differ only in locations in which the same assignment statement is inserted are such that these differences do not impact dynamic data dependencies (according to the semantics in Figure 8), then these two programs either both pass the test or both fail the test. A proof for the above proposition is given in Section C.

Example 3.9 (Test-equivalent programs w.r.t. $\overset{t}{\sim}_{deps}$). Consider a program p_1 defined as

```
x := y;
if y > 0 then
  y := x + 1
else skip fi
```

a program p_2 defined as

```
if y > 0 then
  x := y;
  y := x + 1
else skip fi
```

and a test $t := (\sigma_{in}, \lambda \sigma. \sigma(y) = 3)$ where the input state $\sigma_{in} := \{x \mapsto 1, y \mapsto 2\}$. These programs are test-equivalent w.r.t. the dependency-based test-equivalence relation for the test t , since they differ only in the location of the assignment $x := y$ and the following relation holds:

$$\langle p_1, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle \sigma_{out}, \{l_1, l_2\}, _, _ \rangle$$

where l_1 is the location of $x := y$ in p_1 , l_2 is the location of $x := y$ in p_2 , $\sigma_{out} := \{x \mapsto 2, y \mapsto 3\}$. This relation can be established by the derivation tree in Figure 9.

$$\begin{array}{c}
\frac{\text{Left}(y)}{\langle y, \sigma_1, \{l_1, l_2\}, \{l_1\}, true \rangle \Downarrow_{deps} \langle 2, \{l_1, l_2\}, \{l_1\}, true \rangle} \text{VAR-NLEFT} \quad \frac{\text{Right}(y)}{\langle y := x + 1, \sigma_1, \{l_1, l_2\}, \{l_1\}, true \rangle \Downarrow_{deps} \langle \sigma_{out}, \{l_1, l_2\}, \emptyset, false \rangle} \text{ASSIGN} \\
\frac{\langle y > 0, \sigma_1, \{l_1, l_2\}, \{l_1\}, true \rangle \Downarrow_{deps} \langle true, \{l_1, l_2\}, \{l_1\}, true \rangle} \text{OP} \quad \frac{\langle 0, \dots \rangle \Downarrow 0} \text{NUM} \quad \frac{\langle x + 1, \sigma_1 \rangle \Downarrow 3} \text{ASSIGN-LR-EXE} \\
\frac{\langle \text{if } y > 0 \text{ then } y := x + 1 \text{ else skip } fi, \sigma_1, \{l_1, l_2\}, \{l_1\}, true \rangle \Downarrow_{deps} \langle \sigma_{out}, \{l_1, l_2\}, \emptyset, false \rangle} \text{IF-TRUE} \\
\frac{\text{VAR}}{\langle y, \sigma_{in} \rangle \Downarrow 2} \text{ASSIGN} \\
\frac{\text{Inerted}(x := y)}{\langle x := y, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle \sigma_1, \{l_1, l_2\}, \emptyset, true \rangle} \text{ASSIGN-INS} \quad \frac{\langle \text{if } y > 0 \text{ then } y := x + 1 \text{ else skip } fi, \sigma_1, \{l_1, l_2\}, \{l_1\}, true \rangle \Downarrow_{deps} \langle \sigma_{out}, \{l_1, l_2\}, \emptyset, false \rangle} \text{SEQ} \\
\frac{\langle p_1, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle \sigma_{out}, \{l_1, l_2\}, \emptyset, false \rangle} \text{SEQ}
\end{array}$$

Fig. 9. Derivation tree for dependency-based test-equivalence relation in Example 3.9 ($\sigma_{in} := \{x \mapsto 1, y \mapsto 2\}$, $\sigma_1 := \{x \mapsto 2, y \mapsto 2\}$, $\sigma_{out} := \{x \mapsto 2, y \mapsto 3\}$).

$$\begin{aligned}
M(s_1; s_2) &= \cup_{s' \in M(s_1)} \{ s'; s_2 \} \cup \cup_{s' \in M(s_2)} \{ s_1; s' \} \\
M(\text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}) &= \cup_{s' \in M(s_1)} \{ \text{if } e \text{ then } s' \text{ else } s_2 \text{ fi} \} \cup \cup_{s' \in M(s_2)} \{ \text{if } e \text{ then } s_1 \text{ else } s' \text{ fi} \} \\
M(\text{while } e \text{ do } s \text{ od}) &= \cup_{s' \in M(s)} \{ \text{while } e \text{ do } s' \text{ od} \} \\
M(s \mid s \text{ is not a sequence}) &= M_{\text{EXPRESSION}}(s) \cup M_{\text{REFINEMENT}}(s) \cup M_{\text{GUARD}}(s) \cup M_{\text{ASSIGNMENT}}(s) \\
M_{\text{EXPRESSION}}(v := e) &= \cup_{e' \in \mathcal{E}} \{ v := e' \} \\
M_{\text{EXPRESSION}}(\text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}) &= \cup_{e' \in \mathcal{E}} \{ \text{if } e' \text{ then } s_1 \text{ else } s_2 \text{ fi} \} \\
M_{\text{EXPRESSION}}(\text{while } e \text{ do } s \text{ od}) &= \cup_{e' \in \mathcal{E}} \{ \text{while } e' \text{ do } s \text{ od} \} \\
M_{\text{REFINEMENT}}(\text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}) &= \cup_{e' \in \mathcal{E}} \{ \text{if } e \text{ and } e' \text{ then } s_1 \text{ else } s_2 \text{ fi, if } e \text{ or } e' \text{ then } s_1 \text{ else } s_2 \text{ fi} \} \\
M_{\text{REFINEMENT}}(\text{while } e \text{ do } s \text{ od}) &= \cup_{e' \in \mathcal{E}} \{ \text{while } e \text{ and } e' \text{ do } s \text{ od, while } e \text{ or } e' \text{ do } s \text{ od} \} \\
M_{\text{GUARD}}(v := e) &= \cup_{e' \in \mathcal{E}} \{ \text{if } e' \text{ then } v := e \text{ else skip fi} \} \\
M_{\text{ASSIGNMENT}}(s) &= \cup_{e' \in \mathcal{E}, v' \in \mathcal{V}} \{ v' := e'; s \}
\end{aligned}$$

Fig. 10. Search space definition via transformation schemas M .

3.5 Composing relations

In the proposed test-equivalence analysis framework, several relations can be composed in a mutually-reinforcing fashion. By combining several analyses we can produce a more effective (coarse-grained) partitioning of the space of program modifications into test-equivalence classes.

Definition 3.10 (Composition of test-equivalence relations). Let \mathcal{P} be a finite search space and $\sim_1, \sim_2, \dots, \sim_n$ be test-equivalence relations in \mathcal{P} . A *composition* of $\sim_1, \sim_2, \dots, \sim_n$ is a test-equivalence relation \sim^* such that it is the transitive closure of the union of $\sim_1, \sim_2, \dots, \sim_n$:

$$\sim^* := \left(\bigcup_i \sim_i \right)^*$$

In this work, we define a test-equivalence relation \sim^* as a composition of the relations \sim_{value} and \sim_{deps} introduced in Section 3.3 and Section 3.4:

$$\sim^* := (\sim_{\text{value}} \cup \sim_{\text{deps}})^*$$

An example of applying \sim^* is given in Section 2.3.

The definition above is non-constructive in that it does not define an algorithm of computing test-equivalence classes w.r.t. \sim^* . One possible way to compute test-equivalence classes is to separately evaluate a given program using the semantics in Figure 7 and Figure 8 and merge the results, however this requires multiple program executions. Instead, we use a more efficient approach by defining an augmented semantics of \mathcal{L} for \sim^* by combining the semantics in Figure 7 and Figure 8 as shown in Section B. This semantics enables our approach to compute a test-equivalence class w.r.t. \sim^* via a single execution.

4 PATCH GENERATION

Automated program repair techniques search for patches in spaces of candidate program modifications. A search space in program repair is defined as in the following.

Definition 4.1 (Search space). A *search space* is a finite set of syntactically different programs obtained by applying a given transformation function $M : \mathcal{L} \rightarrow 2^{\mathcal{L}}$ to the buggy program.

ALGORITHM 2: Enumerative patch generation**Input:** search space \mathcal{P} , cost function κ , test-suite T **Output:** ordered set of repairs R

```

1  $R := \emptyset$ ;
2 while  $\mathcal{P} \neq \emptyset$  do
3    $p := \text{pick}(\mathcal{P}, \kappa)$ ;
4    $\text{isPassing} := \text{true}$ ;
5   foreach  $t \in T$  do
6      $\text{isPassing} := \text{eval}(p, t)$ ;
7     if  $\neg \text{isPassing}$  then
8       break;
9   if  $\text{isPassing}$  then
10     $R := R \cup \{p\}$ ;
11     $\mathcal{P} := \mathcal{P} \setminus \{p\}$ ;
12 return  $R$ ;

```

Previous systems (e.g. SPR/Prophet [15, 17] and SemFix/Angelix [22, 25]) defined their search spaces through *parameterized transformation schemas* such that each schema transforms a given program into a program with “holes” and the “holes” are filled with expressions using a program synthesizer. We define our search space in a similar fashion via the function M in Figure 10 (\mathcal{E} indicates the synthesized expressions).

Definition 4.2 (Optimal program repair). Let T be a test-suite (a set of tests), $p \in \mathcal{L}$ be a buggy program ($\exists t \in T. \neg \text{Pass}[p, t]$), $M : \mathcal{L} \rightarrow 2^{\mathcal{L}}$ be a transformation function, $\mathcal{P} := M(p)$ be the corresponding search space, $\kappa : \mathcal{P} \rightarrow \mathbb{R}$ be a cost function. The goal of *optimal program repair* is to find a repair $p' \in \mathcal{P}$ such that $\forall t \in T. \text{Pass}[p', t]$ and $\kappa(p')$ is minimal among all such programs.

Our patch generation algorithm systematically explore the search space by (1) evaluating candidates in the order defined by the prioritization (cost function) starting from the highest priority patch and (2) skipping redundant executions by on-the-fly identifying test-equivalence classes w.r.t. a given test-equivalence relation.

In order to abstract over various optimal synthesis methodologies (synthesis with cost), we assume that there is a function *pick* that for a given set of programs \mathcal{P} and a cost function κ , returns a program from \mathcal{P} with the minimal value of κ .

Consider a baseline enumerative patch generation method described in Algorithm 2. It takes a patch space (Definition 4.1), a cost function and a test-suite as inputs, and outputs a sequence of search space elements that pass all the given tests ordered according to the cost function. First, the algorithm initializes the list of output repairs R . Second, it iterates through the search space by (1) picking the best (the lowest cost according to κ) remaining candidate using *pick* and (2) evaluating the candidate with the tests using *eval*. Finally, it output the list of found plausible patches R .

The overall workflow of our approach is described in Algorithm 3. Our algorithm takes a patch space, a cost function, a test-suite and a test-equivalence relation as inputs and outputs a sequence of search space elements that pass all the given tests ordered according to the cost function. Compared with Algorithm 2, our test-equivalence based algorithm also maintains sets C and \bar{C} for each test. $C(t)$ is a set of test-equivalence classes (therefore, C is a set of sets of programs) in which all candidates pass t ; $\bar{C}(t)$ is the corresponding set of failing test-equivalence classes. First, our algorithm initializes the list of output repairs R and the passing and failing test-equivalence classes C and \bar{C} for all tests t . Second, it iterates through the search space by (1) picking the best (the lowest

ALGORITHM 3: Systematic exploration with partitioning into test-equivalence classes**Input:** search space \mathcal{P} , cost function κ , test-suite T , test-equivalence relation $\overset{\dagger}{\sim}$ **Output:** ordered set of repairs R

```

1  $R := \emptyset$ ;
2 foreach  $t \in T$  do
3    $C(t), \overline{C}(t) := \emptyset, \emptyset$ ;
4 while  $\mathcal{P} \neq \emptyset$  do
5    $p := \text{pick}(\mathcal{P}, \kappa)$ ;
6   if  $\exists t. \forall c \in \overline{C}(t). p \in c$  then
7      $\mathcal{P} := \mathcal{P} \setminus \{p\}$ ;
8     continue;
9   foreach  $t \in T$  do
10    if  $\forall c \in C(t). p \in c$  then
11      continue;
12     $\text{isPassing}, [p] := \widetilde{\text{eval}}(p, t, \mathcal{P}, \overset{\dagger}{\sim})$ ;
13    if  $\text{isPassing}$  then
14       $C(t) := C(t) \cup \{[p]\}$ ;
15    else
16       $\overline{C}(t) := \overline{C}(t) \cup \{[p]\}$ ;
17      break;
18  if  $\forall t. \forall c \in C(t). p \in c$  then
19     $R := R \cup \{p\}$ ;
20   $\mathcal{P} := \mathcal{P} \setminus \{p\}$ ;
21 return  $R$ ;
```

cost according to κ) remaining candidate using *pick* and (2) evaluating the candidate with the tests and computing test-equivalence classes.

For a given candidate, in order to identify the result of a test execution and the corresponding test-equivalence class, the algorithm evaluates the candidate using the function $\widetilde{\text{eval}}$. The function $\widetilde{\text{eval}}$ takes a program p , a test t , a search space \mathcal{P} and a test-equivalence relation $\overset{\dagger}{\sim}$ and returns the result of executing p with t (as a boolean value *isPassing*) and a set of programs $[p]$ such that $[p]$ is a test-equivalence class of p in \mathcal{P} w.r.t. the relation $\overset{\dagger}{\sim}$. The concrete implementation of $\widetilde{\text{eval}}$ depends on the relation $\overset{\dagger}{\sim}$ and is formally described for the relations $\overset{\dagger}{\sim}_{\text{value}}$ and $\overset{\dagger}{\sim}_{\text{deps}}$ in Definition 4.3 and Definition 4.4 respectively.

The test-equivalence classes are used at two steps of search space exploration. First, after a next candidate is picked, the algorithm checks if the candidate belongs to any of the existing failing classes (line 6). If the candidate is in a failing class of at least one test, evaluation of this candidate is omitted. Second, after a next test is selected for evaluating a candidate, the algorithm checks if the candidate is in a passing class of the given test (line 10). If a candidate is in a passing class of a test, then the algorithm omits execution of this candidate with this test. We now discuss the function $\widetilde{\text{eval}}$ specifically for the two notions of test-equivalence we have studied, value-based test-equivalence and dependency-based test-equivalence.

Definition 4.3 (Value-based test-equivalence analysis). Let \mathcal{P} be a search space, $p \in \mathcal{P}$ be a program, $t = (\sigma_{in}, \phi)$ be a test. Let e be an expression in p such that $\exists p' \in \mathcal{P} \exists e' \in \text{Expr}. p' = p[e/e']$. Then,

value-based test-equivalence analysis \widetilde{eval} is defined as follows:

$$\begin{aligned} \widetilde{eval}(p, t, \mathcal{P}, \overset{\dagger}{\sim}_{value}) &= (\phi(\sigma_{out}), \bigcup_{e' \in \mathcal{C}} \{p[e/e']\}), \text{ given that} \\ \text{Modified} &:= \lambda x. x = e, \\ \mathcal{E} &:= \{e' \mid \exists p' \in \mathcal{P}. p' = p[e/e']\}, \\ \langle p, \sigma_{in}, \mathcal{E} \rangle &\Downarrow_{value} \langle \sigma_{out}, C \rangle \end{aligned}$$

In this analysis, we identify a test-equivalence class of a program with an expression e in the space of all programs in \mathcal{P} that differ only in e . The test-equivalence class is computed by passing the set of all “alternative” expressions \mathcal{E} as an argument of \Downarrow_{value} . Note that in this definition we explicitly select an expression e , substitution of which are analyzed for test-equivalence. For each element of a search space produced by the transformation function M in Figure 10, there is always at most one such e . We now discuss the function \widetilde{eval} for dependency-based test-equivalence.

Definition 4.4 (Dependency-based test-equivalence analysis). Let \mathcal{P} be a search space, $p \in \mathcal{P}$ be a program, $t = (\sigma_{in}, \phi)$ be a test. Let p' be a program, l_1 be a location such that $p = p'[l_1/v := e; l_1]$ and $\exists p'' \in \mathcal{P}. \exists l_2 \in p''. p'' = p'[l_2/v := e; l_2]$. Then, dependency-based test-equivalence analysis \widetilde{eval} is:

$$\begin{aligned} \widetilde{eval}(p, t, \mathcal{P}, \overset{\dagger}{\sim}_{deps}) &= (\phi(\sigma_{out}), \bigcup_{l \in L'} \{p'[l/v := e; l]\}), \text{ given that} \\ \text{Inserted} &:= \lambda s. s = “v := e”, \text{ Left} := \lambda v'. v' = v, \text{ Right} := \lambda v'. v' \in \text{Var}(e), \\ L &:= \{l \mid l \in p' \wedge \exists p'' \in \mathcal{P}. p'' = p'[l/v := e; l]\} \\ \langle p, \sigma_{in}, L, \emptyset, \text{false} \rangle &\Downarrow_{deps} \langle \sigma_{out}, L', _, _ \rangle \end{aligned}$$

In this analysis, we identify a test-equivalence class of a program with an assignment $v := e$ in the space of all programs in \mathcal{P} that differ only in locations of this assignment. The test-equivalence class is computed by passing the set of all “alternative” locations L as an argument of \Downarrow_{deps} .

Finally, note that Algorithm 3 can be used in different ways. The output of the algorithm is a sequence of plausible patches R ordered according to the function κ . A sequence of repairs can be used to provide several patch suggestions for developers. The number of suggested repairs can be controlled by introducing a limit and breaking from the main loop when the required number of plausible patches is found. Certain applications may require generation of all plausible patches (e.g. in order to narrow candidates through test generation [31]). In this case, the algorithm can be modified so that it outputs whole test-passing partitions instead of single patches.

5 IMPLEMENTATION

We have implemented the described approach in a tool called `f1x` (pronounced as [ɛf-waɪn-ɛks]) for the C programming language.

Analysis. Our implementation of the proposed test-equivalence analyses is built upon a combination of static (source code) and dynamic instrumentation. Specifically, to implement the augmented semantics in Section 3.3 for the relation $\overset{\dagger}{\sim}_{value}$, we apply the transformation schemas M (Figure 10) to the source code of the buggy program and replace “holes” with calls to a procedure implementing the value-projection operator. To implement the augmented semantics in Section 3.4 for the relation $\overset{\dagger}{\sim}_{deps}$, we implemented a dynamic instrumentation using Pin [18] that tracks reads and writes of the variables involved in assignment synthesis.

Search space. The goal of this work was to design and evaluate test-equivalence relations for transformations used in existing program repair systems. Our system combines the transformation

schemas of SPR/Prophet and Angelix (we studied implementation of these systems in order to closely reproduce their search spaces), described as follows.

EXPRESSION Modify an existing side-effect free integer expression or condition (adopted from Angelix). A variant of $M_{EXPRESSION}$ in Figure 10 for C programs. Partitioned into test-equivalence classes based on the expression values using $\overset{t}{\sim}_{value}$.

REFINEMENT Append a disjunct/conjunct to an existing condition (adopted from Prophet). A variant of $M_{REFINEMENT}$ in Figure 10 for C programs. Partitioned into test-equivalence classes based on the condition values using $\overset{t}{\sim}_{value}$.

GUARD Add an if-guard for an existing statement (adopted from Angelix and Prophet). A variant of M_{GUARD} in Figure 10 for C programs. Partitioned into test-equivalence classes based on the condition values using $\overset{t}{\sim}_{value}$.

ASSIGNMENT Insert an assignment statement (adopted from Prophet⁷). A variant of $M_{ASSIGNMENT}$ in Figure 10 for C programs. Partitioned into test-equivalence classes using $\overset{t}{\sim}^*$.

INITIALIZATION Insert memory initialization (adopted from Prophet). Not partitioned.

FUNCTION Replace a function call with another function (adopted from Prophet). Not partitioned.

The two last transformation schemas adopted from Prophet are not partitioned by our algorithm, since they generate relatively small search spaces. Our transformations differ from that of SPR/Prophet in the following ways: (1) Prophet implements a transformation schema for inserting guarded return statements. Although our algorithm can partition these transformations using the relation $\overset{t}{\sim}_{value}$, such transformations were shown to frequently generate overfitting patches [38] and therefore we exclude them from our search space; (2) Prophet implements a transformation that copies existing program statements. Since such statements can be arbitrarily complex and they cannot be partitioned by our algorithm, we do not include this transformation.

Cost function. Several techniques have been proposed to increase the probability of generating correct repairs by prioritizing patches [3, 17, 21]. For our system we implement an approach that assigns higher priority to smaller changes [21]:

$$\kappa(p) := \text{distance}(p, p_{orig})$$

where p is a patched program (an element of the search space), p_{orig} is the original program, *distance* is defined as the number of added, modified and deleted AST nodes.

6 EXPERIMENTAL EVALUATION

We evaluate our approach in terms of the following research questions:

- (RQ1) What are the effectiveness and efficiency of our approach compared with state-of-the-art program repair systems?
- (RQ2) Does our approach scale to larger search spaces compared with state-of-the-art systems? Does test-equivalence relation enable higher scalability of our implementation?
- (RQ3) What is the impact of each test-equivalence relation on the number of test executions performed by our algorithm?

6.1 Evaluation setup

Our evaluation compares f1x against three repair approaches: Angelix, Prophet and GenProg-AE. These repair techniques are chosen as they use different repair algorithms including symbolic analysis (Angelix), machine-learning (Prophet) and genetic algorithm (GenProg). We evaluate all

⁷Prophet generates new assignments by copying and modifying existing assignments. Instead, f1x *synthesizes* assignments and therefore its search space includes a superset of assignments that can be generated by Prophet.

Table 1. Subject programs and their basic statistics

Program	Description	LOC	Defects	Tests	Execution cost (sec)
libtiff	Image processing library	77K	24	78	8.18
lighttpd	Web server	62K	9	295	29.75
php	Interpreter	1046K	44	8471	427.25
gmp	Math library	145K	2	146	53.52
gzip	Data compression utility	491K	5	12	0.43
python	Interpreter	407K	11	35	156.46
wireshark	Network packet analyzer	2814K	7	63	9.92
fcc	Compiler	97K	3	773	240.27

repair approaches on the GenProg ICSE'12 benchmark [12] for our evaluation because it includes defects from large real-world projects, and was designed for systematic evaluation of program repair tools. Moreover, the test suites in this benchmark were independently augmented to prevent repair tools from generating implausible patches [28]. The benchmark consists of 105 defects from eight subjects (i.e. libtiff, lighttpd, PHP, gmp, gzip, python, wireshark, and fbc) which have developer-written test suites. Table 1 shows the statistics of each evaluated subject. The column "Execution cost" denotes the time taken to execute the test-suite for a given subject.

We selected the following systems and their configurations for evaluation:

- F1X** f1x that implements search with test-equivalence partitioning described in Algorithm 3.
- F1X^E** f1x^E is a variant of f1x that enumerates changes without test-equivalence partitioning (Algorithm 2). It is considered to evaluate implementation-independent effect of partitioning.
- ANG** Angelix 1.1 [22] that implements a symbolic path exploration and prioritizes syntactically small changes.
- PR** Prophet 0.1 [17] that implements value search (a variant of path exploration) for conditional expressions and patch prioritization based on machine learning.
- PR*** Prophet* that is a variant of Prophet that disables transformations for (1) inserting overfitting return insertions and (2) copying complex statements except for assignments. This variant is considered to match the transformation implemented in F1X/F1X^E, since the search space of F1X/F1X^E is effectively the combination of the search spaces of PR* and ANG.
- GP** GenProg-AE 3.0 [42] that implements a group of analysis techniques to avoid evaluating functionally-equivalent patches (as opposite to test-equivalent as in our approach). Compared to the earlier version of GenProg that uses genetic algorithm [14] that is inherently stochastic, GenProg-AE leverages a deterministic repair algorithm.

We run all the configurations (F1X, F1X^E, ANG, PR, PR*, GP) in two modes:

- Stop-after-first-found** The algorithm terminates after finding the first patch. This mode represents the usual program repair usage scenario.
- Full exploration** The algorithm terminates after searching through the entire search space. This mode allows us to obtain data that is independent on (1) the exploration order and (2) whether a plausible patch is present in the search space.

We reuse the configurations from previous studies for running Angelix, Prophet and GenProg-AE [28, 42]. As Prophet takes a correctness model as input to prioritizes patches akin to the provided model, we used the default model that is publicly available⁸.

We conduct all experiments on Intel[®] Xeon[™] CPU E5-2660 machines running Ubuntu 14.04, and use a 10 hours timeout for running each configuration.

⁸Prophet website: <http://rhino.csail.mit.edu/prophet-rep/>

Table 2. Effectiveness of program repair approaches.

Subject	Plausible						Equivalent to human					
	F1X	F1X ^E	ANG	PR	PR*	GP	F1X	F1X ^E	ANG	PR	PR*	GP
libtiff	13	10	10	5	3	5	5	3	3	2	1	0
lighttpd	5	3	-	4	4	4	0	0	-	0	0	0
php	15	7	10	18	15	7	6	3	4	10	6	2
gmp	2	1	2	2	2	1	2	1	2	1	1	0
gzip	3	2	2	2	2	2	2	0	1	1	1	0
python	5	1	-	6	5	3	0	0	-	0	0	1
wireshark	4	4	4	4	4	4	0	0	0	0	0	0
fbc	1	1	-	1	1	1	1	1	-	1	1	0
Overall	49	29	28	42	36	27	16	8	10	15	10	3

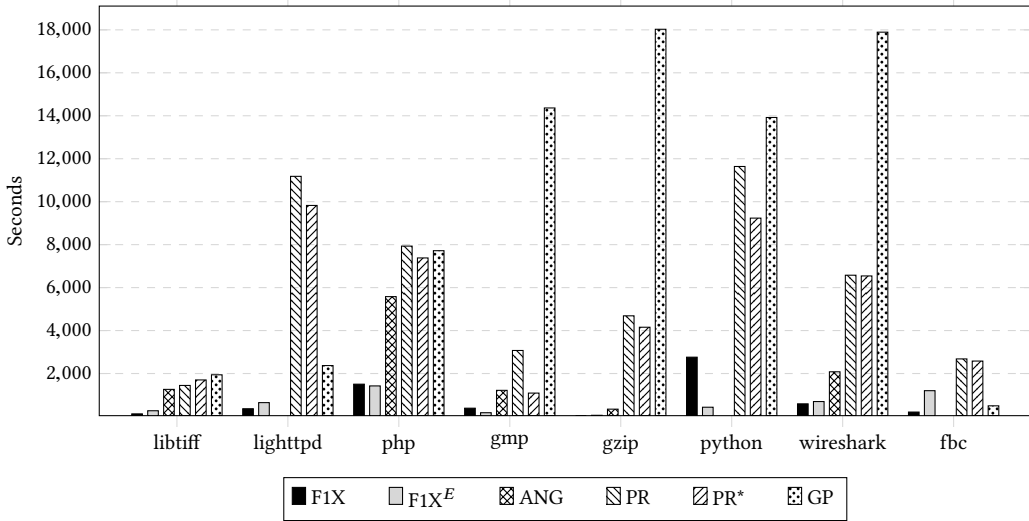


Fig. 11. Average patch generation time.

6.2 Effectiveness and efficiency (RQ1)

Table 2 summarizes the effectiveness results for F1X, F1X^E, ANG, PR, PR* and GP executed in the stop-after-first-found mode. The second through seventh columns denote the number of plausible patches generated by each repair approach, while the eighth through thirteenth columns represent the number of patches syntactically equivalent to the human patches. As Angelix does not support lighttpd, python and fbc, the corresponding cells for these subjects are marked with “-”. The overall results illustrate that F1X generates the highest number of plausible patches compared to all other evaluated repair approaches. The “Equivalent to human” column in table 2 shows that F1X generates 8 more human-like patches than F1X^E, 6 more human-like patches than ANG, 1 more human-like patch than PR, 2 more human-like patches than PR* and 13 more human-like patches than GP.

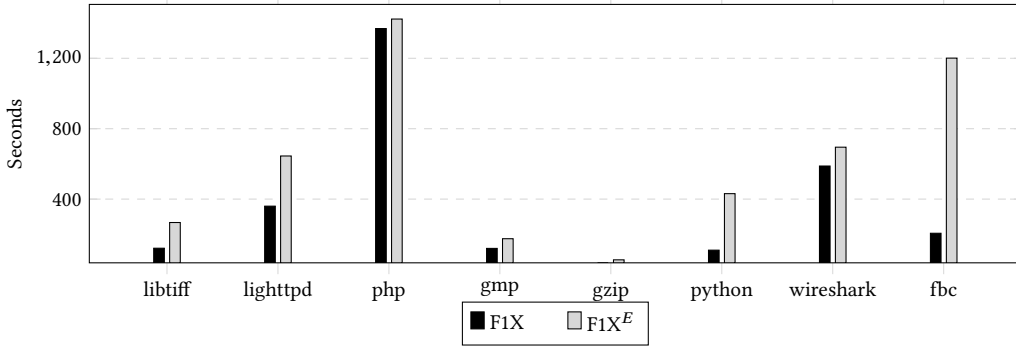


Fig. 12. Average time for patches generated by both F1X and F1X^E.

We attribute the high number of patches generated by F1X to the larger patch space supported by F1X compared to other approaches. Since F1X combines the search spaces of ANG and PR*, it fixes all defects that are fixed by either of these tools. Note that F1X finds more patches than F1X^E within the time limit due to the performance gain from our partitioning.

Figure 11 illustrates the average patch generation time for the configurations. The x-axis of Figure 11 represents the eight subjects in the benchmark, while the y-axis shows the average time taken to generate a patch for all defects for a given subject where each bar depicts a patch generation approach. Overall, the average patch generation time for F1X is significantly shorter than all other repair approaches. For instance, F1X requires only 121 seconds on average to generate a patch for libtiff, while ANG takes 1262 seconds (F1X is $\frac{1262}{121}=10.5X$ faster than ANG). Meanwhile, PR* takes 1701 seconds on average to produce a patch for libtiff (F1X is $\frac{1701}{121}=14X$ faster than PR*). Notably, F1X is 16X faster than GP for libtiff (GP takes 1940 seconds on average to generate a patch for libtiff). The average patch generation time for PR is slightly higher compared to PR* as it searches through a slightly larger patch space.

The results shown in Figure 11 validate our claim that F1X is able to achieve significant improvement on the patch generation time due to its efficient search algorithm. F1X and F1X^E demonstrate a comparable average time of patch generation. Note that for some subjects (e.g. python), the average time of F1X^E is lower than that of F1X. This is because F1X^E finds a subset of patches found by F1X, and patches found exclusively by F1X contribute to its higher average time. However, when considering only patches found by both F1X and F1X^E, F1X shows consistently better performance, as shown in Figure 12.

RQ1: Compared with existing automated program repair tools, f1x generates more patches since its search space is a combination the search spaces of ANG and PR*. Despite a larger search space, it finds patches in an order of magnitude faster.

6.3 Exploration speed (RQ2)

Definition 6.1 (Explored candidates). We say that a candidate patch is *explored* if the algorithm identified whether the patch passes all given tests or fails at least one. Note that we only consider candidate patches in which the source code modification is executed by all given failing tests.

Table 3 shows the exploration statistics for F1X, F1X^E, PR, PR* and GP (we exclude ANG because the search space for Angelix is encoded via logical constraints). The second through sixth columns depicts the data for the stop-after-first-found mode, while the seventh through eleventh columns

Table 3. Exploration statistics of program repair tools in stop-after-first-found and full exploration modes.

Subject	Stop-after-first-found mode					Full exploration mode				
	F1X	F1X ^E	PR	PR*	GP	F1X	F1X ^E	PR	PR*	GP
libtrif	$\frac{44675}{243}=183.8$	$\frac{1241}{1272}=1.0$	$\frac{14980}{2946}=5.1$	$\frac{6995}{4625}=1.5$	$\frac{1371}{1302}=1.1$	$\frac{400786}{759}=528$	$\frac{23872}{24512}=1.0$	$\frac{62521}{23763}=2.6$	$\frac{31403}{16777}=1.9$	$\frac{4564}{4657}=1.0$
lighttpd	$\frac{1085}{220}=4.9$	$\frac{512}{616}=0.8$	$\frac{4645}{1295}=3.6$	$\frac{6269}{2563}=2.4$	$\frac{359}{383}=0.9$	$\frac{112667}{850}=132$	$\frac{10283}{10580}=1.0$	$\frac{45338}{8377}=5.4$	$\frac{29106}{10126}=2.9$	$\frac{2337}{2371}=1.0$
php	$\frac{2843}{1717}=1.7$	$\frac{3407}{5184}=0.7$	$\frac{1186}{48070}=0.0$	$\frac{2566}{56492}=0.0$	$\frac{88}{7736}=0.0$	$\frac{77472}{193}=401$	$\frac{9801}{12139}=0.8$	$\frac{6378}{58671}=0.1$	$\frac{4942}{73898}=0.1$	$\frac{2057}{3497}=0.6$
gmp	$\frac{5517}{173}=31.9$	$\frac{30}{31}=1.0$	$\frac{3215}{2002}=1.6$	$\frac{2064}{1260}=1.6$	$\frac{9128}{9172}=1.0$	$\frac{80448}{1140}=70$	$\frac{16934}{23994}=0.7$	$\frac{17526}{14523}=1.2$	$\frac{11423}{10156}=1.1$	$\frac{2841}{2841}=1.0$
gzip	$\frac{13071}{123}=106.3$	$\frac{241}{251}=1.0$	$\frac{5892}{1340}=4.4$	$\frac{5284}{890}=5.9$	$\frac{7605}{5735}=1.3$	$\frac{518803}{568}=913$	$\frac{18501}{24789}=0.7$	$\frac{59467}{28716}=2.1$	$\frac{48491}{40574}=1.2$	$\frac{9007}{9098}=1.0$
python	$\frac{940}{231}=4.1$	$\frac{29}{37}=0.8$	$\frac{4211}{7389}=0.6$	$\frac{6928}{8771}=0.8$	$\frac{3327}{3319}=1.0$	$\frac{105842}{1617}=65$	$\frac{8381}{12494}=0.7$	$\frac{13367}{9213}=1.5$	$\frac{9928}{9050}=1.1$	$\frac{7823}{7824}=1.0$
wshark	$\frac{828}{139}=6.0$	$\frac{595}{597}=1.0$	$\frac{3744}{2460}=1.5$	$\frac{3948}{2434}=1.6$	$\frac{4253}{4297}=1.0$	$\frac{689925}{620}=1112$	$\frac{14094}{14003}=1.0$	$\frac{23043}{13099}=1.8$	$\frac{18554}{17442}=1.1$	$\frac{2008}{2213}=0.9$
fbc	$\frac{948}{434}=2.1$	$\frac{650}{891}=0.7$	$\frac{1111}{672}=1.65$	$\frac{524}{394}=1.33$	$\frac{22}{788}=0.0$	$\frac{50195}{1312}=38$	$\frac{20195}{24521}=0.8$	$\frac{892}{589}=1.51$	$\frac{766}{499}=1.54$	$\frac{852}{852}=1.0$

represents the data for the full exploration mode. Each cell in the second through eleventh columns is of the form $\frac{X}{Y}=Z$ where X represents the average number of “explored candidates” by a repair approach, Y represents the average number of test executions performed by a repair approach, Z denotes the exploration speed (computed by the ratio of the number of “explored candidates” over the number of test executions). The average for the stop-after-first-found is computed among the fixed defects, whereas the average for the full exploration mode is computed among all defects.

In general, F1X has an order of magnitude higher exploration speed compared to all other patch generation approaches in both the stop-after-first-found mode and the full exploration mode. For example, in the full exploration mode, F1X requires on average only 620 test executions to explore 689925 candidates for wireshark. For the same subject, PR requires 13099 test executions for exploring 23043 candidates, PR* requires 17442 test executions for exploring 18554 candidates and GP requires 2213 test executions for exploring 2008 candidates. The efficiency of exploration can also be indirectly shown by comparing the average number of plausible patches found in full exploration mode. For this, we compared the results of F1X and PR, since these configuration found the largest number of patches. F1X generate 2265 plausible patches on average, while PR generate only 9 plausible patches on average.

To enable generation of more patches, an ideal repair approach should scale to larger spaces by exploring more candidates within the time budget. The data in Table 3 shows that F1X scales to larger search spaces, since it explores more candidates within the time limit due to fewer number of test executions. This explains the effectiveness and the efficiency of F1X compared with other tools shown in Table 2 and Figure 11. Recall that F1X^E is a variant of F1X without test-equivalence partitioning. For the same search space, F1X^E explores less candidates than F1X within the time limit since it requires more test executions. From this observation, we conclude that test-equivalence partitioning is responsible for the higher scalability of F1X.

Table 4. Effect of equivalence relations on the number of test executions on Libtiff

Transformation	Relation	Locations	Candidates	Partitions	Test executions
$M_{EXPRESSION}$	$\overset{t}{\sim}_{value}$	109	428641	424	739
$M_{REFINEMENT}$	$\overset{t}{\sim}_{value}$	45	48942	87	154
M_{GUARD}	$\overset{t}{\sim}_{value}$	87	106389	199	347
$M_{ASSIGNMENT}$	$\overset{t}{\sim}_{value}$	121	17181	3436	5983
	$\overset{t}{\sim}_{deps}$	121	17181	1809	2308
	$\overset{t}{\sim}^*$	121	17181	480	671

RQ2: f1x scales to larger search space by exploring more candidates with fewer number of test executions.

6.4 Effect of equivalent relation (RQ3)

To investigate the effect of each equivalence relation on the number of test executions, we conduct another experiment on all 24 defects in Libtiff. We selected Libtiff, because due to the structure of its source code, Libtiff subjects have the largest total number of candidates patches in the search space ($9618864 = 400786 * 24$), where 400786 is the average number of candidates per version, 24 is the number of versions (see Table 3). Our goal is to determine if there is a single equivalence relation that dominantly contributes to the reduction in the number test executions or if a composition of these relations will result in greater reduction.

Table 4 shows the effect of the three equivalent relations ($\overset{t}{\sim}_{value}$, $\overset{t}{\sim}_{deps}$, and $\overset{t}{\sim}^*$) in F1X on the average number of test executions for Libtiff. For each transformation from Figure 10, the table demonstrates the number of locations in which the transformation was applied (the “Locations” column), the test-equivalence relations that was applied for the search space produced by this transformation (the “Relation” column), the number of different candidate patches generated (the “Candidates” column), the number of partitions that were identified for the failing test (the “Partitions” column) and the number of tests required to explore all the corresponding candidates with the whole test suite (the “Test executions” column).

These results demonstrate that for the transformations $M_{EXPRESSION}$, $M_{REFINEMENT}$ and M_{GUARD} our algorithm produces a small number of test-equivalence classes for the failing tests (with 500-1000 elements in each partition on average) which also resulted in a small number of executions for the whole test suite. For the relation $M_{ASSIGNMENT}$, the reduction is less significant, however the composition of relations $\overset{t}{\sim}^*$ is significantly more efficient (the number of test executions is significantly less) than the individual relations $\overset{t}{\sim}_{value}$, $\overset{t}{\sim}_{deps}$.

RQ3: f1x partitions a large number of patches into a small number of test-equivalence classes using two relations: $\overset{t}{\sim}_{value}$ and $\overset{t}{\sim}_{deps}$. The composition of these two relations is more effective than the individual relations.

7 RELATED WORK

Program synthesis. Existing program synthesis techniques can be used to generate patches by directly searching in patch spaces, however this approach has limitations as explained in the following. First, the cost of test execution in a typical program repair problem is substantially higher

than that in program synthesis, since the subjects on which program repair is carried out today are significantly larger (e.g. a single test execution for PHP interpreter from GenProg benchmark [12] takes 1-10 seconds on commodity hardware, while a typical solution in SyGuS-Comp competition can be executed in 10^{-6} seconds). Second, the complexity of repaired programs makes it infeasible to apply precise deductive techniques. For instance, Sketch [33] fills “holes” in sketches (partial programs), and can be potentially applied to generate repairs for identified suspicious statements. However, it translates programs into boolean formulas and therefore can repair only relatively small programs [5]. Since program synthesis algorithms may not be directly applicable to program repair, they are used as parts of program repair algorithms for filling “holes” in programs based on inferred specification [15, 22]. In our technique, we do not use program synthesis as a black box, but integrate synthesis with program analysis by imposing additional requirements on the underlying synthesizer: support for the value-projection operator (Section 3.2). Since our implementation uses an enumerative program synthesis, it is straightforward to realize such operators. However, other techniques can also be used for this purpose. For instance, FlashMeta [26] compactly represents its search space as version space algebra (VSA) [23]. Moreover, it defines the operation *Filter* over this representation that is effectively the value-projection operator, therefore it can be potentially used in our algorithm as a more efficient representation of the space of program modifications.

Program repair search algorithms. Syntax-based techniques generate patches by enumerating and testing syntactic changes. Since (1) repair tools have to explore large search spaces to address many classes of defects and (2) test execution has high cost for large real-world programs, they scale to relatively small search spaces. GenProg-AE [42] eliminates redundant executions by identifying functionally-equivalent patches via lightweight analyses. Instead of functional equivalence, our technique applies test-equivalence, which is a weaker and therefore a more effective relation (produces larger equivalence classes). *Semantics-based* techniques split search into two phases. First, they localize suspicious statements and *infer specification* for the identified statements that captures the property of “passing the test suite”. Such specification can be expressed as logical constraints [25] or *angelic values* [15, 44]. Second, they apply off-the-shelf program synthesizers in order to modify the selected statements according to this specification. To infer specification, existing techniques perform path exploration by altering test executions. Since the number of execution paths in programs can be infinite, these methods are subject to the *path explosion problem*. For instance, Nopol [44], SPR [15] and Prophet [17] enumerate values of conditional expressions (which is a special case of path exploration) in order to find angelic values that enable the program to pass the failing test. As shown in Section 2.1, such techniques may perform a large number of redundant executions that are avoided by our algorithm. SemFix [25] and Angelix [22] are semantics-based techniques relying on symbolic execution and SMT-based synthesis. Since our methodology does not use symbolic methods, it is orthogonal to SemFix and Angelix from that point of view of underlying analysis. Existing syntax and semantic-based techniques are limited to modifying side-effect free expressions. For instance, they can only generate new assignments by copying them from other parts of the program. Meanwhile, we demonstrate that test-equivalence can scale assignment synthesis using a combination of value- and dependency-based analyses.

Program repair prioritization approaches. In order to address the test over-fitting problem [32], various techniques have been proposed to prioritize patches that are more likely to be correct. For instance, DirectFix [21] prioritizes candidate patches based on syntactic distance, Qlose [3] prioritizes patches bases on semantic distance, Prophet [17] utilizes information learned from human patches, ACS [43] prioritizes patches based on information mined from previous version and API documentation, S3 [11] prioritizes patches based on a combination of syntactic and semantic properties. Meanwhile, relifix [35] leverages different program version for generating

patches, whereas SEMGRAFT [20] uses references implementation as implicit correctness criteria. Droix [36] fixes crashes in Android apps using Android activity lifecycle management rules for patch prioritization. Our technique finds the best patch (the global optimum) in its search space according to an arbitrary cost function. Meanwhile, previous techniques applied to large real-world programs did not provide such guarantees. Techniques based on genetic programming [2, 14] and random search [27] guarantee only a local optimum by definition. Semantics-based repair techniques may miss the global optimum, which is shown in Section 2.2.

Program repair using types and formal specification. Several approaches utilize temporal logic formulas [7], contracts [41] and types [29] to guide program repair. Our test-equivalence analysis can potentially optimize these approaches. Besides, our test-driven patch generation algorithm can be used in a counterexample-guided refinement loop [1] to repair programs based on given formal specification.

Mutation testing. The scalability of program repair is related to scalability of mutation testing, since mutation testing also evaluates a large number of program modifications. To address this problem in mutation testing, a common approach is to reduce the number of mutation operators (transformations) to avoid redundant executions [24]. This approach may not be suitable for program repair, because program repair search spaces have to be rich enough to enable generation of non-trivial human-like repairs. Test-equivalence has been applied to scale mutation testing. Mutant analysis by Just et al. [8] performs a pre-pass that partitions mutants based on infected states, which can be thought of as a variant of \sim_{value} relation. More recent techniques [19, 40] extend this approach by performing more fine-grained partitioning. Our method adapts value-based test-equivalence relation that has been used in mutation testing to program repair by integrating it with syntax-guided program synthesis through value-projector operator (Section 3.2). Apart from that, we introduce a new dependency-based test-equivalence relation for partitioning more complex program modification that might be required to fix real bugs.

Compiler testing. Equivalence modulo inputs (EMI) [9, 10] (a variant of test-equivalence) was successfully applied to compiler testing. Specifically, for a compiler, a program and an input, it generates input-equivalent variants of the program by altering unexecuted statements and checks that these variants compiled by the compiler produce the same outputs. A recent extension of this approach [34] modifies executed statements by inserting random code guarded by a condition that is evaluated into *false* in the context of the given test, which can be considered as an application of value-based test-equivalence relation. Our dependency-based test-equivalence relation and the proposed composition of several analyses might be used to increase the effectiveness of compiler testing by synthesizing non-trivial input-equivalent program modifications.

8 DISCUSSION AND FUTURE WORK

We envision the following design of a future *general-purpose program repair system* (a system that is able to address many kinds of defects in commodity software). This system will (1) implement a large number of transformations to address many kinds of defects, (2) implement test-equivalence analyses for these transformations to ensure scalability and (3) implement intelligent search space prioritization strategies over the patch space, to address the overfitting problem. This work is a step towards such a design. In our future works, we plan to investigate the following aspects.

Test-equivalence relations. The effectiveness of the relation \sim_{value} (the size of test-equivalence classes it induces) may depend on the size of the output domain of modified expressions. Since it identifies equivalence based on concrete values, it may not be effective for expressions of large

output domains (e.g. strings). In future, we plan to investigate a generalization of this relation that defines two changes to be test-equivalent iff they drive test execution along the same path. Such a relation will be computed using dynamic symbolic execution [4].

Transformations. Existing program repair systems provide very limited support for repairing function calls. We hypothesize that a generalization of the relation $\overset{\dagger}{\sim}_{deps}$ may enable repair systems to extend search spaces by incorporating function call transformations in a scalable fashion.

Current limitations. Although we demonstrated that our approach based on test-equivalence significantly outperforms previous repair techniques, it has several limitations. The proposed analysis assumes deterministic test execution. Besides, our algorithm is designed to search in finite spaces of candidate patches. It can be potentially generalized for infinite spaces by encoding them symbolically. Finally, the current algorithm is designed to synthesize single-line patches (involving a single modification). It can be potentially extended for multi-line modifications using the approach of Angelix [22].

Subject Programs. The subjects for the evaluation (GenProg ICSE'12 benchmark [12]) were previously used for evaluating related approaches [15, 17, 22]. GenProg ICSE'12 benchmark was constructed to address generalizability concerns in evaluation of repair tools [12]. Moreover, the test suites were independently augmented [28] to avoid generation of implausible patches. Nevertheless, a possible threat to validity is that our results may not generalize for other programs and defects. In future, it may be worthwhile to evaluate our approaches on other relevant benchmarks [13, 37]

9 CONCLUSION

Traditionally, many problems in software testing which involve exploring huge search spaces, such as mutation testing, have been shown to benefit from test equivalence [8, 9]. In this paper, we have adopted the notion of test-equivalence for program repair, and provided several test-equivalence relations: based on runtime values, based on dynamic data dependencies and their composition. We also proposed an algorithm of automatic patch generation based on on-the-fly test-equivalence analysis. This enables us to explore larger patch spaces in substantially less time, thereby significantly enhancing the practicality of modern day program repair technology.

ACKNOWLEDGEMENTS

The authors would like to thank Julia Lawall for her comments on earlier drafts of this paper. This work is supported in part by Office of Naval Research grant ONRG-NICOP-N62909-18-1-2052. The authors would also like to thank Airbus Helicopters for the sponsorship of this work as part of Airbus Project Skyways.

REFERENCES

- [1] Rajeev Alur, Rastislav Bodik, Garvit Juniwal, Milo MK Martin, Mukund Raghothaman, Sanjit A Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. 2013. Syntax-guided synthesis. In *FMCAD*. IEEE, 1–8.
- [2] Andrea Arcuri and Xin Yao. 2008. A novel co-evolutionary approach to automatic software bug fixing. In *CEC*. IEEE, 162–168.
- [3] Loris D'Antoni, Roopsha Samanta, and Rishabh Singh. 2016. Qlose: program repair with quantitative objectives. In *CAV*. Springer, 383–401.
- [4] Patrice Godefroid, Nils Klarlund, and Koushik Sen. 2005. DART: directed automated random testing. In *PLDI*. ACM, 213–223.
- [5] Jinru Hua and Sarfraz Khurshid. 2016. A sketching-based approach for debugging using test cases. In *ATVA*. Springer, 463–478.
- [6] Susmit Jha, Sumit Gulwani, Sanjit A Seshia, and Ashish Tiwari. 2010. Oracle-guided component-based program synthesis. In *ICSE*. 215–224.

- [7] Barbara Jobstmann, Andreas Griesmayer, and Roderick Bloem. 2005. Program repair as a game. In *CAV*. Springer, 226–238.
- [8] René Just, Michael D Ernst, and Gordon Fraser. 2014. Efficient mutation analysis by propagating and partitioning infected execution states. In *ISSTA*. ACM, 315–326.
- [9] Vu Le, Mehrdad Afshari, and Zhendong Su. 2014. Compiler validation via equivalence modulo inputs. In *PLDI*. ACM, 216–226.
- [10] Vu Le, Chengnian Sun, and Zhendong Su. 2015. Finding deep compiler bugs via guided stochastic program mutation. In *OOPSLA*. ACM, 386–399.
- [11] Xuan Bach Dinh Le, Duc Hiep Chu, David Lo, Claire Le Goues, and Willem Visser. 2017. S3: syntax- and semantic-guided repair synthesis via programming by example. In *FSE*. ACM.
- [12] Claire Le Goues, Michael Dewey-Vogt, Stephanie Forrest, and Westley Weimer. 2012. A systematic study of automated program repair: Fixing 55 out of 105 bugs for 8 each. In *ICSE*. IEEE, 3–13.
- [13] Claire Le Goues, Neal Holtschulte, Edward K Smith, Yuriy Brun, Premkumar Devanbu, Stephanie Forrest, and Westley Weimer. 2015. The ManyBugs and IntroClass benchmarks for automated repair of C programs. *TSE* 41, 12 (2015), 1236–1256.
- [14] Claire Le Goues, ThanhVu Nguyen, Stephanie Forrest, and Westley Weimer. 2012. Genprog: a generic method for automatic software repair. *TSE* 38, 1 (2012), 54–72.
- [15] Fan Long and Martin Rinard. 2015. Staged program repair with condition synthesis. In *FSE*. ACM, 166–178.
- [16] Fan Long and Martin Rinard. 2016. An analysis of the search spaces for generate and validate patch generation systems. In *ICSE*. ACM, 702–713.
- [17] Fan Long and Martin Rinard. 2016. Automatic patch generation by learning correct code. In *POPL*. ACM, 298–312.
- [18] Chi-Keung Luk, Robert Cohn, Robert Muth, Harish Patil, Artur Klauser, Geoff Lowney, Steven Wallace, Vijay Janapa Reddi, and Kim Hazelwood. 2005. Pin: building customized program analysis tools with dynamic instrumentation. In *PLDI*. ACM, 190–200.
- [19] Yu-Seung Ma and Sang-Woon Kim. 2016. Mutation testing cost reduction by clustering overlapped mutants. *Journal of Systems and Software* 115 (2016), 18–30.
- [20] Sergey Mechtaev, Manh-Dung Nguyen, Lars Noller, Yannic Grunske, and Abhik Roychoudhury. 2018. Semantic Program Repair Using a Reference Implementation. In *ICSE*.
- [21] Sergey Mechtaev, Jooyong Yi, and Abhik Roychoudhury. 2015. Directfix: looking for simple program repairs. In *ICSE*. IEEE, 448–458.
- [22] Sergey Mechtaev, Jooyong Yi, and Abhik Roychoudhury. 2016. Angelix: Scalable Multiline Program Patch Synthesis via Symbolic Analysis. In *ICSE*.
- [23] Tom M Mitchell. 1982. Generalization as search. *Artificial intelligence* 18, 2 (1982), 203–226.
- [24] Elfurjani S Mresa and Leonardo Bottaci. 1999. Efficiency of mutation operators and selective mutation strategies: An empirical study. *STVR* 9, 4 (1999), 205–232.
- [25] Hoang Duong Thien Nguyen, Dawei Qi, Abhik Roychoudhury, and Satish Chandra. 2013. SemFix: program repair via semantic analysis. In *ICSE*. 772–781.
- [26] Oleksandr Polozov and Sumit Gulwani. 2015. FlashMeta: a framework for inductive program synthesis. *OOPSLA* (2015), 107–126.
- [27] Yuhua Qi, Xiaoguang Mao, Yan Lei, Ziyang Dai, and Chengsong Wang. 2014. The strength of random search on automated program repair. In *ICSE*. ACM, 254–265.
- [28] Zichao Qi, Fan Long, Sara Achour, and Martin Rinard. 2015. An analysis of patch plausibility and correctness for generate-and-validate patch generation systems. In *ISSTA*. ACM, 24–36.
- [29] Alex Reinking and Ruzica Piskac. 2015. A type-directed approach to program repair. In *CAV (1)*. 511–517.
- [30] Reudismam Rolim, Gustavo Soares, Loris D’Antoni, Oleksandr Polozov, Sumit Gulwani, Rohit Gheyi, Ryo Suzuki, and Björn Hartmann. 2017. Learning syntactic program transformations from examples. In *ICSE*. IEEE Press, 404–415.
- [31] David Shriver, Sebastian Elbaum, and Kathryn T Stolee. 2017. At the end of synthesis: narrowing program candidates. In *ICSE NIER*. IEEE Press, 19–22.
- [32] Edward K Smith, Earl T Barr, Claire Le Goues, and Yuriy Brun. 2015. Is the cure worse than the disease? overfitting in automated program repair. In *FSE*. ACM, 532–543.
- [33] Armando Solar-Lezama. 2008. *Program synthesis by sketching*. University of California, Berkeley.
- [34] Chengnian Sun, Vu Le, and Zhendong Su. 2016. Finding compiler bugs via live code mutation. In *OOPSLA*. ACM, 849–863.
- [35] Shin Hwei Tan and Abhik Roychoudhury. 2015. relifix: Automated repair of software regressions. In *ICSE*. IEEE, 471–482.
- [36] Shin Hwei Tan, Abhik Roychoudhury, Zhen Dong, and Xiang Gao. 2018. Repairing Crashes in Android Apps. In *ICSE*.

- [37] Shin Hwei Tan, Jooyong Yi, Sergey Mechtaev, Abhik Roychoudhury, et al. 2017. Codeflaws: a programming competition benchmark for evaluating automated program repair tools. In *ICSE Companion*. 180–182.
- [38] Shin Hwei Tan, Hiroaki Yoshida, Mukul R Prasad, and Abhik Roychoudhury. 2016. Anti-patterns in search-based program repair. In *FSE*. ACM, 727–738.
- [39] Yida Tao, Jindae Kim, Sunghun Kim, and Chang Xu. 2014. Automatically generated patches as debugging aids: a human study. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. ACM, 64–74.
- [40] Bo Wang, Yingfei Xiong, Yangqingwei Shi, Lu Zhang, and Dan Hao. 2017. Faster mutation analysis via equivalence modulo states. *ISSTA (2017)*.
- [41] Yi Wei, Yu Pei, Carlo A Furia, Lucas S Silva, Stefan Buchholz, Bertrand Meyer, and Andreas Zeller. 2010. Automated fixing of programs with contracts. In *ISSTA*. ACM, 61–72.
- [42] Westley Weimer, Zachary P Fry, and Stephanie Forrest. 2013. Leveraging program equivalence for adaptive program repair: Models and first results. In *ASE*. IEEE, 356–366.
- [43] Yingfei Xiong, Jie Wang, Runfa Yan, Jiachen Zhang, Shi Han, Gang Huang, and Lu Zhang. 2017. Precise condition synthesis for program repair. In *ICSE*. IEEE Press, 416–426.
- [44] Jifeng Xuan, Matias Martinez, Favio Demarco, Maxime Clément, Sebastian Lamelas, Thomas Durieux, Daniel Le Berre, and Martin Monperrus. 2016. Nopol: Automatic Repair of Conditional Statement Bugs in Java Programs. *TSE (2016)*.
- [45] Jooyong Yi, Umair Ahmed, Amey Karkare, Shin Hwei Tan, and Abhik Roychoudhury. 2017. A feasibility study of using automated program repair for introductory programming assignments. In *FSE*. ACM.

A ALTERNATIVE IMPLEMENTATIONS OF VALUE-PROJECTION OPERATOR

We define the value-projection operator $\Pi_{\sigma,n}^{value}$ (Section 3.2) for SMT-based synthesis [6] and DSL-based synthesis [26].

A.1 Symbolic synthesis

In SMT-based component-based synthesis [6], the search space is represented implicitly through logical constraints. To synthesize an expression from a given specification, it solves the following formula:

$$\exists e \in \mathcal{E}. \bigwedge_{\sigma, n \in \text{spec}} \langle e, \sigma \rangle \Downarrow n$$

To support quantification over expressions, it uses *location variables* (denoted via the prefix l) to encode all expressions constructed from given components C (e.g. $+$, $-$, variables, etc) as follows.

$$\begin{aligned} \phi_{\text{wpf}} &:= \phi_{\text{range}} \wedge \phi_{\text{cons}} \wedge \phi_{\text{acyc}} \\ \phi_{\text{range}} &:= \bigwedge_{c \in C} \left(0 \leq lc^{\text{out}} < |C| \wedge \bigwedge_{k \in [1, NI(c)]} 0 \leq lc_k^{\text{in}} < |C| \right) \\ \phi_{\text{cons}} &:= \bigwedge_{(c,s) \in C \times C, c \neq s} lc^{\text{out}} \neq ls^{\text{out}} \\ \phi_{\text{acyc}} &:= \bigwedge_{c \in C, k \in [1, NI(c)]} lc^{\text{out}} > lc_k^{\text{in}} \\ \phi_{\text{conn}} &:= \bigwedge_{\substack{(c,s) \in C \times C \\ k \in [1..NI(s)]}} lc^{\text{out}} = ls_k^{\text{in}} \Rightarrow c^{\text{out}} = s_k^{\text{in}} \end{aligned}$$

The algorithm also imposes library constraint (ϕ_{lib}) that capture semantics of given components. For example, for the component $c := h_1 + h_2$, the library constraint is $c^{\text{out}} = c_1^{\text{in}} + c_1^{\text{in}}$. Given the above constraints, the search space is encoded through the formula ϕ defined as $\phi := \phi_{\text{wpf}} \wedge \phi_{\text{lib}} \wedge \phi_{\text{conn}}$.

In this context, the value-projection operator can be implemented as follows:

$$\Pi_{\sigma,n}^{value}(\phi) := \phi \wedge e^{out} = n \wedge \bigwedge_{\text{variable } v \in C} v^{out} = \sigma(v)$$

where e^{out} captures the output of the synthesized expression. Effectively, this operator conjoins the formula representing the search space with the input-output relation represented via σ and n .

A.2 DSL-based synthesis

In FlashMeta [26] synthesis framework, the search space is compactly represented via version space algebra (VSA) [23]. Assume that the set of expression \mathcal{E} is defined through applications of operators to a given set of variables; we denote an operator as F .

The grammar for a version space algebra \tilde{N} is defined as

$$\tilde{N} := \{ e_1, \dots, e_k \} \mid U(\tilde{N}_1, \dots, \tilde{N}_k) \mid F_{\triangleright\triangleleft}(\tilde{N}_1, \dots, \tilde{N}_k)$$

such that

- $e \in \{ e_1, \dots, e_k \}$ if $\exists i. e = e_i$;
- $e \in U(\tilde{N}_1, \dots, \tilde{N}_k)$ if $\exists i. e \in \tilde{N}_i$;
- $e \in F_{\triangleright\triangleleft}(\tilde{N}_1, \dots, \tilde{N}_k)$ if $e = F(e_1, \dots, e_n) \wedge \forall i. e_i \in \tilde{N}_i$.

For version space algebra \tilde{N} , FlashMeta implements a clustering operator denoted as $\tilde{N}|_{\sigma}$. $\tilde{N}|_{\sigma}$ is a mapping from values to version space algebras $\{ n_1 \mapsto \tilde{N}_1, \dots, n_k \mapsto \tilde{N}_k \}$ such that

- $\tilde{N} = \tilde{N}_1 \cup \dots \cup \tilde{N}_k$;
- $\tilde{N}_i \cap \tilde{N}_j = \emptyset$ for all $i \neq j$;
- $\forall e \in \tilde{N}_i. \langle e, \sigma \rangle \Downarrow n_i$;
- $\forall i, j. i \neq j \rightarrow n_i \neq n_j$.

In this context, the value-projection operator can be implemented as follows:

$$\Pi_{\sigma,n}^{value}(\tilde{N}) := \tilde{N}|_{\sigma}(n)$$

B SEMANTICS FOR COMPOSITION OF TEST-EQUIVALENCE RELATIONS

Section 3.5 presents a non-constructive definition of the composition of test-equivalence relations $\overset{\dagger}{\sim}^*$ (it does not define an algorithm for computing test-equivalence partitions). Although it is possible to compute test-equivalence partitions w.r.t. $\overset{\dagger}{\sim}^*$ by separately applying the value-based test-equivalence analysis in Figure 7 and dependency-based test-equivalence analysis in Figure 8, this approach would require performing multiple program executions. A more efficient approach is to compute test-equivalence partitions for $\overset{\dagger}{\sim}^*$ directly thorough a combination of the semantics in Figure 7 and Figure 8 as shown in Figure 13. For a given program p with an assignment $v := \varepsilon$, this semantics identifies a mapping $\{ l_1 \mapsto \mathcal{E}_1, \dots, l_k \mapsto \mathcal{E}_k \}$ that denotes all test-equivalent insertions of assignments of $v := e'$ at each location l_i so that $e' \in \mathcal{E}_i$. In this semantics, l and c denote such mappings, $l \cap c$ denotes $\lambda x. l(x) \cup c(x)$ and $l \setminus c$ denotes $\lambda x. l(x) \setminus c(x)$.

For this semantics of $\overset{\dagger}{\sim}^*$, test-equivalence analysis for the patch generation algorithm in Section 4 can be defined as follows:

Definition B.1 (Composed test-equivalence analysis). Let \mathcal{P} be a search space, $p \in \mathcal{P}$ be a program, $t = (\sigma_{in}, \phi)$ be a test. Let p' be a program, l_1 be a location such that $p = p'[l_1/v := \varepsilon; l_1]$ and $\exists p'' \in$

<p>VAR-LEFT-EXE</p> $\frac{Left(v) \wedge x}{\langle v, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma(v), l \cap c, \emptyset, false \rangle}$ <p>OP</p> $\frac{\langle e_1, \sigma, l, c, x \rangle \Downarrow_* \langle n_1, l_1, c_1, x_1 \rangle \quad \langle e_2, \sigma, l_1, c_1, x_1 \rangle \Downarrow_* \langle n_1, l_1, c_1, x_1 \rangle \quad n_3 = n_1 \text{ op } n_2}{\langle e_1 \text{ op } e_2, \sigma, l, c, x \rangle \Downarrow_* \langle n_3, l_2, c_2, x_2 \rangle}$	<p>VAR-LEFT-NEXE</p> $\frac{Left(v) \wedge \neg x}{\langle v, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma(v), l \setminus c, \emptyset, x \rangle}$ <p>ASSIGN-LR-EXE</p> $\frac{-Inserted(v := e) \wedge (Left(v) \vee Right(v)) \wedge x \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma', l \cap c, \emptyset, false \rangle}$ <p>ASSIGN-LR-NEXE</p> $\frac{-Inserted(v := e) \wedge (Left(v) \vee Right(v)) \wedge \neg x \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma', l \setminus c, \emptyset, false \rangle}$ <p>ASSIGN-NLR</p> $\frac{-Inserted(v := e) \wedge \neg Left(v) \wedge \neg Right(v) \quad \langle e, \sigma, l, c, x \rangle \Downarrow_* \langle n, l_1, c_1, x_1 \rangle}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma[v \mapsto n], l_1, c_1, x_1 \rangle}$	<p>VAR-NLEFT</p> $\frac{\neg Left(v)}{\langle v, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma(v), l, c, x \rangle}$ <p>SKIP</p> $\frac{}{\langle skip, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma, l, c, x \rangle}$
<p>ASSIGN-INS</p> $\frac{Inserted(v := e) \quad \langle v := e, \sigma \rangle \Downarrow \sigma'}{\langle v := e, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma', l, c, true \rangle}$	<p>SEQ</p> $\frac{\langle s_1, \sigma, l, c, x \rangle \Downarrow_{value*} \langle \sigma_1, l_1, c_1, x_1 \rangle \quad \langle s_2, \sigma_1, l_1, c_1, x_1 \rangle \Downarrow_{value*} \langle \sigma_2, l_2, c_2, x_2 \rangle}{\langle s_1 ; s_2, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma_2, l_2, c_2, x_2 \rangle}$ <p>IF-TRUE</p> $\frac{\langle e, \sigma, l, c, x \rangle \Downarrow_* \langle true, l_1, c_1, x_1 \rangle \quad \langle s_1, \sigma, l_1, c_1, x_1 \rangle \Downarrow_{value*} \langle \sigma_1, l_2, c_2, x_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, c \rangle \Downarrow_* \langle \sigma_1, l_2, c_2, x_2 \rangle}$ <p>IF-FALSE</p> $\frac{\langle e, \sigma, l, c, x \rangle \Downarrow_* \langle false, l_1, c_1, x_1 \rangle \quad \langle s_2, \sigma, l_1, c_1, x_1 \rangle \Downarrow_{value*} \langle \sigma_2, l_2, c_2, x_2 \rangle}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma_2, l_2, c_2, x_2 \rangle}$	<p>WHILE-TRUE</p> $\frac{\langle e, \sigma, l, c, x \rangle \Downarrow_* \langle true, l_1, c_1, x_1 \rangle \quad \langle s_1, \sigma, l_1, c_1, x_1 \rangle \Downarrow_{value*} \langle \sigma_1, l_2, c_2, x_2 \rangle \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma_1, l_2, c_2, x_2 \rangle \Downarrow_{value*} \langle \sigma_2, l_3, c_3, x_3 \rangle}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma_2, l_3, c_3, x_3 \rangle}$ <p>WHILE-FALSE</p> $\frac{\langle e, \sigma, l, c, x \rangle \Downarrow_* \langle false, l', c', x' \rangle \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma, l', c', x' \rangle}{\langle \text{while } e \text{ do } s \text{ od}, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma, l', c', x' \rangle}$ <p>NUM</p> $\frac{}{\langle n, \sigma \rangle \Downarrow n}$ <p>STMT-L</p> $\frac{\langle s, \sigma, l, c, x \rangle \Downarrow_* \langle \sigma', l', c', x' \rangle \quad c'' = \lambda x. \text{if } x = s \text{ then } c'(x) \cap \Pi_{\sigma, n}^{value}(\epsilon) \text{ else } c'(x)}{\langle s, \sigma, l, c, x \rangle \Downarrow_{value*} \langle \sigma', l', c'', x' \rangle}$

Fig. 13. Augmented semantics of \mathcal{L} for computing test-equivalence classes w.r.t. $\stackrel{t}{\sim}_*$. v – variables, n – integer values, e – expressions, s – statements, σ – program states, l, c – functions from locations to sets of expressions, x – boolean values, $Inserted$ – predicate over statements, $Left, Right$ – predicates over variables, ϵ – right-hand side of inserted assignment.

$\mathcal{P}. \exists l_2 \in p''. \exists e'. p'' = p'[l_2/v := e':l_2]$. Then, dependency-based test-equivalence analysis \widetilde{eval} is:

$$\widetilde{eval}(p, t, \mathcal{P}, \stackrel{t}{\sim}_{deps}) = (\phi(\sigma_{out}), \bigcup_{l \in p} \bigcup_{e \in L'(l)} \{p'[l/v := e:l]\}), \text{ given that}$$

$$\begin{aligned} Inserted &:= \lambda s. s = "v := e", \quad Left := \lambda v'. v' = v, \quad Right := \lambda v'. v' \in Var(e), \\ L &:= \lambda l. \{e \mid \exists p'' \in \mathcal{P}. p'' = p'[l/v := e:l]\} \\ \langle p, \sigma_{in}, L, \lambda x. \emptyset, false \rangle &\Downarrow_* \langle \sigma_{out}, L', _, _ \rangle \end{aligned}$$

In this analysis, we identify a test-equivalence class of a program with an assignment $v := e'$ in the space of all programs in \mathcal{P} that differ only in the expressions e' and the locations of this assignment. The test-equivalence class is computed by passing the “alternative” mapping from locations to expressions L as an argument of \Downarrow_* .

$$\begin{aligned}
& \text{Trace}\left(\frac{\langle v, \sigma \rangle \Downarrow \sigma(v)}{\text{VAR}^i}\right) := \{ \lambda \sigma. \sigma[x_i \mapsto \sigma(v)] \} & \text{Trace}\left(\frac{\langle n, \sigma \rangle \Downarrow n}{\text{NUM}^i}\right) := \{ \lambda \sigma. \sigma[x_i \mapsto n] \} \\
& \text{Trace}\left(\frac{\langle e_1, \sigma \rangle \Downarrow n_1 \quad \langle e_2, \sigma \rangle \Downarrow n_2 \quad n_3 = n_1 \text{ op } n_2}{\langle e_1 \text{ op } e_2, \sigma \rangle \Downarrow n_3} \text{OP}^i\right) := \text{Trace}(\langle e_1, \sigma \rangle \Downarrow n_1) @ \text{Trace}(\langle e_2, \sigma \rangle \Downarrow n_2) @ \{ \lambda \sigma. \sigma[x_i \mapsto \sigma(x_i) \text{ op } \sigma(x_r)] \} \\
& \text{Trace}\left(\frac{\langle \text{skip}, \sigma \rangle \Downarrow \sigma}{\text{SKIP}^i}\right) := \{ \lambda \sigma. \sigma \} & \text{Trace}\left(\frac{\langle e, \sigma \rangle \Downarrow n}{\langle v := e, \sigma \rangle \Downarrow \sigma[v \mapsto n]} \text{ASSIGN}^i\right) := \text{Trace}(\langle e, \sigma \rangle \Downarrow n) @ \{ \lambda \sigma. \sigma[v \mapsto n] \} \\
& \text{Trace}\left(\frac{\langle s_1, \sigma \rangle \Downarrow \sigma_1 \quad \langle s_2, \sigma_1 \rangle \Downarrow \sigma_2}{\langle s_1 ; s_2, \sigma \rangle \Downarrow \sigma_2} \text{SEQ}^i\right) := \text{Trace}(\langle s_1, \sigma \rangle \Downarrow \sigma_1) @ \text{Trace}(\langle s_2, \sigma_1 \rangle \Downarrow \sigma_2) \\
& \text{Trace}\left(\frac{\langle e, \sigma \rangle \Downarrow \text{true} \quad \langle s_1, \sigma \rangle \Downarrow \sigma_1}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \Downarrow \sigma_1} \text{IF-TRUE}^i\right) := \text{Trace}(\langle e, \sigma \rangle \Downarrow \text{true}) @ \text{Trace}(\langle s_1, \sigma \rangle \Downarrow \sigma_1) \\
& \text{Trace}\left(\frac{\langle e, \sigma \rangle \Downarrow \text{false} \quad \langle s_2, \sigma \rangle \Downarrow \sigma_2}{\langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \Downarrow \sigma_2} \text{IF-FALSE}^i\right) := \text{Trace}(\langle e, \sigma \rangle \Downarrow \text{false}) @ \text{Trace}(\langle s_2, \sigma \rangle \Downarrow \sigma_2) \\
& \text{Trace}\left(\frac{\langle e, \sigma \rangle \Downarrow \text{true} \quad \langle s_1, \sigma \rangle \Downarrow \sigma_1 \quad \langle \text{while } e \text{ do } s \text{ od}, \sigma_1 \rangle \Downarrow \sigma_2}{\langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \Downarrow \sigma_2} \text{WHILE-TRUE}^i\right) := \text{Trace}(\langle e, \sigma \rangle \Downarrow \text{true}) @ \text{Trace}(\langle s_1, \sigma \rangle \Downarrow \sigma_1) @ \text{Trace}(\langle \text{while } e \text{ do } s \text{ od}, \sigma_1 \rangle \Downarrow \sigma_2) \\
& \text{Trace}\left(\frac{\langle e, \sigma \rangle \Downarrow \text{false}}{\langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \Downarrow \sigma} \text{WHILE-FALSE}^i\right) := \text{Trace}(\langle e, \sigma \rangle \Downarrow \text{false})
\end{aligned}$$

Fig. 14. Definition of *Trace*. @ — list concatenation operator, v — variables, n — integer values, e — expressions, s — statements, σ — program states, i — index of current rule application, l, r — indexes of child rules application (left and right), x_i — variable introduced for each rule application, $\text{Trace}(\langle e_1, \sigma \rangle \Downarrow n_1)$ indicates an application of *Trace* to the derivation tree for $\langle e_1, \sigma \rangle \Downarrow n_1$.

C PROPERTIES OF TEST-EQUIVALENCE SEMANTICS

We introduce a linear representation of a derivation tree that we refer to as a *trace*. A trace of a derivation tree D (for the semantics in Figure 6) is denoted as $\text{Trace}(D)$ (Figure 14). The trace of the derivation tree in Example 3.2 is a sequence of functions

$$\{ f_{\text{VAR}}^0, f_{\text{NUM}}^1, f_{\text{OP}}^2, f_{\text{ASSIGN}}^3, f_{\text{VAR}}^4, f_{\text{ASSIGN}}^5, f_{\text{SEQ}}^6 \}$$

where f_{RULE}^i is a *state transformation function* for the rule RULE with the index i in the derivation tree computed in depth-first order from left to right. Note that $\sigma_{\text{out}} = f_{\text{VAR}}^0 \circ f_{\text{NUM}}^1 \circ f_{\text{OP}}^2 \circ f_{\text{ASSIGN}}^3 \circ f_{\text{VAR}}^4 \circ f_{\text{ASSIGN}}^5 \circ f_{\text{SEQ}}^6(\sigma_{\text{in}})$ (by construction).

C.1 Value-based test-equivalence relation

This section describes properties and their proofs for the semantics of value-based test-equivalence relation in Figure 7. In an analogous manner to the previous section, we define the function *Trace* for the value-based semantics in such a way that it transforms each rule application into a function $f_{\text{RULE}}^i : \Sigma \times \mathcal{E} \rightarrow \Sigma \times \mathcal{E}$, where the second element of the tuple denotes the computed class of expressions.

LEMMA C.1. *Let p be a program, e be an expression, $t := (\sigma_{\text{in}}, \phi)$ be a test, $\text{Modified} := \lambda x. x = e$. If $\langle p, \sigma_{\text{in}}, C_{\text{in}} \rangle \Downarrow_{\text{value}} \langle _ , C_{\text{out}} \rangle$ (with the corresponding derivation tree D) and $e \in C_{\text{out}}$, then for each applications of the rule *EXPR-MOD* in the derivation tree D with premises $\langle e, \sigma \rangle \Downarrow n$ and $c' = \Pi_{\sigma, n}^{\text{value}}(c)$, the expression e is in the set c' .*

PROOF. The rule *EXPR-MOD* is either encountered in D or not. If it is not applied in the derivation tree, then $C_{\text{out}} = C_{\text{in}}$ since *EXPR_MOD* is the only rule that affects c . If it is applied in the derivation tree, represent the derivation tree as a trace:

$$\{ f_{\text{R}_0}^0, f_{\text{R}_1}^1, \dots, f_{\text{R}_k}^k \}$$

Then, $(\sigma_{out}, C_{out}) = f_{R_0}^0 \circ f_{R_1}^1 \circ \dots \circ f_{R_k}^k(\sigma_{in}, C_{in})$. Therefore, the set C_{out} can be represented as $\Pi_{\sigma_1, n_1}^{value} \circ \Pi_{\sigma_2, n_2}^{value} \circ \dots \circ \Pi_{\sigma_m, n_m}^{value}(C_{in})$, where each application of the value-projector operator corresponds to an application of the rule EXPR-MOD. Since the value-projector operator outputs a subset of its argument set, $\forall i. e \in \Pi_{\sigma_i, n_i}^{value} \circ \dots \circ \Pi_{\sigma_m, n_m}^{value}(C_{in})$. Therefore, $e \in c'$ for each application of the rule EXPR-MOD in D . \square

LEMMA C.2. *Let p be a program, e be an expression, $t := (\sigma_{in}, \phi)$ be a test, $Modified := \lambda x. x = e$. If $\langle p, \sigma_{in}, C_{in}^1 \rangle \Downarrow_{value} \langle _, C_{out}^1 \rangle$ and $\langle p, \sigma_{in}, C_{in}^2 \rangle \Downarrow_{value} \langle _, C_{out}^2 \rangle$ and $C_{in}^1 \subseteq C_{in}^2$, then $C_{out}^1 \subseteq C_{out}^2$.*

PROOF. It is trivial that for any n and σ if $C_{in}^1 \subseteq C_{in}^2$, then $\Pi_{\sigma, n}^{value}(C_{in}^1) \subseteq \Pi_{\sigma, n}^{value}(C_{in}^2)$. Then, the lemma follows from the fact that C_{out} can be represented as $\Pi_{\sigma_1, n_1}^{value} \circ \Pi_{\sigma_2, n_2}^{value} \circ \dots \circ \Pi_{\sigma_k, n_k}^{value}(C_{in})$. \square

LEMMA C.3. *Let p_1 and p_2 be programs such that $\exists e, e' \in Expr. p_2 = p_1[e'/e]$, $t := (\sigma_{in}, \phi)$ be a test. If $p_1 \overset{t}{\sim}_{value} p_2$ and $\langle p_1, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ (for $Modified := \lambda x. x = e$), then $\langle p_2, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ (for $Modified := \lambda x. x = e'$).*

PROOF. $\langle p_1, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ for $Modified := \lambda x. x = e$ with the corresponding derivation tree D . Therefore, in each application of the rule EXPR-MOD of the derivation tree, $c' = \{e, e'\}$ (Lemma C.1). Consequently, in each occurrence of EXPR-MOD, $\langle e, \sigma \rangle \Downarrow n \Leftrightarrow \langle e', \sigma \rangle \Downarrow n$. Assume that D' is obtained by replacing all occurrences of e with e' in the derivation tree D . Since $\langle e, \sigma \rangle \Downarrow n \Leftrightarrow \langle e', \sigma \rangle \Downarrow n$ in each application of the rule EXPR-MOD, D' is a valid derivation tree for $\langle p_2, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ with $Modified := \lambda x. x = e'$. \square

LEMMA C.4. *Let p be program, $t := (\sigma_{in}, \phi)$ be a test. If $\langle p, \sigma_{in}, _ \rangle \Downarrow_{value} \langle \sigma_{out}, _ \rangle$, then $\langle p, \sigma_{in} \rangle \Downarrow \sigma_{out}$.*

PROOF. A derivation tree for the semantics in Figure 7 can be transformed into a valid derivation tree of the standard semantics (Figure 6) by removing the argument c and changing the applications of EXPR-MOD and EXPR-NMOD into corresponding applications of EXPR. \square

PROPOSITION C.5. *Let \mathcal{P} be a set of programs such that for any $p_1, p_2 \in \mathcal{P}$, $\exists e, e' \in Expr$ such that $p_2 = p_1[e'/e]$, $t := (\sigma_{in}, \phi)$ be a test. Then, $\overset{t}{\sim}_{value}$ as in Definition 3.4 is an equivalence relation in \mathcal{P} .*

PROOF. **Reflexivity.** Let p be a program, e be an expression in p , $Modified := \lambda x. x = e$. Since $\forall e. \langle e, \sigma \rangle \Downarrow n \Leftrightarrow \{e\} = \Pi_{\sigma, n}^{value}(\{e\})$, $c' = \{e\}$ in each application of EXPR-MOD for $c = \{e\}$. Therefore, $\langle p_1, \sigma_{in}, \{e\} \rangle \Downarrow_{value} \langle _, \{e\} \rangle$. Consequently, $p \overset{t}{\sim}_{value} p$. **Symmetry.** Let p_1, p_2 be programs and e_1, e_2 be expressions such that $p_2 = p_1[e_1/e_2]$ and $p_1 \overset{t}{\sim}_{value} p_2$. Assume $\langle p_1, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$. Then, $\langle p_2, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ (Lemma C.3). Therefore, $p_2 \overset{t}{\sim}_{value} p_1$. **Transitivity.** Let p_1, p_2, p_3 be programs and e_1, e_2, e_3 be expressions such that $p_2 = p_1[e_1/e_2]$ and $p_3 = p_2[e_2/e_3]$. Assume $p_1 \overset{t}{\sim}_{value} p_2$ and $p_2 \overset{t}{\sim}_{value} p_3$. Note that for any σ , if $\langle e_1, \sigma \rangle \Downarrow n \wedge c' = \Pi_{\sigma, n}^{value}(\{e_1, e_2, e_3\})$, then $e_2 \in c' \wedge e_3 \in c'$ (Lemma C.1 and Lemma C.2). Thus, $\langle e_1, \sigma \rangle \Downarrow n \Leftrightarrow \{e_1, e_2, e_3\} = \Pi_{\sigma, n}^{value}(\{e_1, e_2, e_3\})$. Therefore, $\forall \sigma. \langle e_2, \sigma \rangle \Downarrow n \Leftrightarrow \{e_2, e_3\} = \Pi_{\sigma, n}^{value}(\{e_2, e_3\})$. Consequently, $\langle p_2, \sigma_{in}, \{e_2, e_3\} \rangle \Downarrow_{value} \langle _, \{e_2, e_3\} \rangle$ for $Modified := \lambda x. x = e_2$ and $p_2 \overset{t}{\sim}_{value} p_3$. \square

Proposition C.5 can be trivially generalized for spaces involving modifications of multiple expressions since the union of several equivalence relations is an equivalence relation.

PROPOSITION C.6. *Let p_1 and p_2 be programs such that $\exists e, e' \in Expr. p_2 = p_1[e'/e]$, $t := (\sigma_{in}, \phi)$ be a test. If $p_1 \overset{t}{\sim}_{value} p_2$, then p_1 and p_2 either both pass t or both fail t .*

PROOF. Assume that $\langle p_1, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ for $Modified := \lambda x. x = e$. Therefore, $\langle p_2, \sigma_{in}, \{e, e'\} \rangle \Downarrow_{value} \langle \sigma_{out}, \{e, e'\} \rangle$ for $Modified := \lambda x. x = e'$ (Lemma C.3), $\langle p_1, \sigma_{in} \rangle \Downarrow \sigma_{out}$ (Lemma C.4) and $\langle p_2, \sigma_{in} \rangle \Downarrow \sigma_{out}$ (Lemma C.4). Since p_1 and p_2 produce the same output state for σ_{in} , they either both pass t or both fail t . \square

PROPOSITION C.7. *The relation $\overset{t}{\sim}_{value}$ is a test-equivalence relation according to Definition 1.1.*

PROOF. Proposition C.5 states that $\overset{t}{\sim}_{value}$ is an equivalence relation and Proposition C.6 states that any two programs equivalent w.r.t. $\overset{t}{\sim}_{value}$ produce indistinguishable results for the given test. Therefore, $\overset{t}{\sim}_{value}$ is a test-equivalence relation. \square

C.2 Dependency-based test-equivalence relation

This section describes properties and their proofs for the semantics of dependency-based test-equivalence relation in Figure 8. In an analogous manner to the previous section, we define the function *Trace* for the dependency-based semantics in such a way that it transforms each rule application into a function $f_{RULE}^i : \Sigma \times Stmt \times Stmt \times \mathbb{B} \rightarrow \Sigma \times Stmt \times Stmt \times \mathbb{B}$, where the seconds and the third elements of the tuple denote global and current test-equivalence classes, and the last boolean element indicates if the inserted assignment has been executed.

LEMMA C.8. *Let p be a program, e be an expression and v be a variable, such that the assignment “ $v := e$ ” is at the location s of p , $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$. Then, $\langle p, \sigma_{in}, \{s\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{s\}, _, _ \rangle$.*

PROOF. Consider a trace of the derivation tree of the relation $\langle p, \sigma_{in}, \{s\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, _, _, _ \rangle$:

$$\{ f_{R_0}, f_{R_1}, \dots, f_{R_k} \}$$

For each i , if $(_, l, c, x) = f_{R_i} \circ \dots \circ f_{R_k}(\sigma_{in}, \{s\}, \emptyset, false)$, then $x = true \Leftrightarrow s \in c$ since x is set to *true* only when the rule ASSIGN-INS is applied (that adds s to c), and x is set to *false* in the rules VAR-LEFT-EXE, VAR-LEFT-NEXE, ASSIGN-LR-EXE and ASSIGN-LR-NEXE (that make $c = \emptyset$). Therefore, when the rules VAR-LEFT-NEXE and ASSIGN-LR-NEXE are applied, $s \notin c$. Thus, s is never removed from l in the derivation tree. Therefore, $\langle p, \sigma_{in}, \{s\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{s\}, _, _ \rangle$. \square

Assume that $\{ f_{R_0}, f_{R_1}, \dots, f_{R_k} \}$ is a trace of some program p in a derivation tree for the semantics in Figure 8. We call a sequence of rule applications $I := \{ f_{R_1}, \dots, f_{R_r} \}$ a *test-equivalent interval* iff

- it does not include the rules VAR-LEFT-EXE, VAR-LEFT-NEXE, ASSIGN-LR-EXE and ASSIGN-LR-NEXE;
- it contains an application of the rule ASSIGN-INS;
- it is maximal (not a sub-interval of another test-equivalent interval).

Note that each application of the rule ASSIGN-INS is included into a test-equivalent interval (possibly consisting only of this rule). We enumerate all test-equivalent intervals in the trace as I_0, I_1, \dots, I_k from left to right.

LEMMA C.9. *Let p_1 and p_2 be programs, e be an expression and v be a variable, such that there is program p with locations l_1, l_2 so that $p_1 = p[l_1/v := e; l_1]$, $p_2 = p[l_2/v := e; l_2]$, $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$, $p_1 \overset{t}{\sim}_{deps} p_2$. Assume that $I_0^1, I_1^1, \dots, I_k^1$ are the test-equivalent intervals in the trace of p_1 for a derivation tree $\langle p_1, \sigma_{in}, _, _ \rangle \Downarrow_{deps} \langle _, _, _, _ \rangle$. Then, for each I_i^1 there is a corresponding test-equivalent interval I_i^2 in the trace of $\langle p_2, \sigma_{in}, _, _ \rangle \Downarrow_{deps} \langle _, _, _, _ \rangle$ such the first rules in the intervals receive the same inputs state for the same program location and the last rules in the intervals produce the same state for the same program location in the trace.*

PROOF. Assume $I_i^1 := \{f_{R_0}^1, \dots, f_{R_k}^1\}$ and $I_i^2 := \{f_{R_0}^2, \dots, f_{R_l}^2\}$. We prove that for any state σ , if $(\sigma', _, _) = f_{R_0}^1 \circ \dots \circ f_{R_k}^1(\sigma, _, _)$ and the rules $f_{R_0}^1$ and $f_{R_0}^2$ are applied for the same program location, then $(\sigma', _, _) = f_{R_0}^2 \circ \dots \circ f_{R_l}^2(\sigma, _, _)$ and $f_{R_l}^1$ and $f_{R_l}^2$ are applied for the same program location. This fact is proved by induction on the rule index in the test-equivalent interval, excluding the applications of ASSIGN-INS. Then, this fact can be used to prove the lemma by induction on the index of the test-equivalent interval. \square

LEMMA C.10. *Let p be program, $t := (\sigma_{in}, \phi)$ be a test. If $\langle p, \sigma_{in}, _, _ \rangle \Downarrow_{deps} \langle \sigma_{out}, _, _, _ \rangle$, then $\langle p, \sigma_{in} \rangle \Downarrow \sigma_{out}$.*

PROOF. A derivation tree for the semantics in Figure 8 can be transformed into a valid derivation tree of the standard semantics (Figure 6) by removing the arguments l , c and x , changing the applications of VAR-LEFT-EXE, VAR-LEFT-NEXE and VAR-NLEFT into corresponding applications of VAR and the applications of ASSIGN-INS, ASSIGN-LR-EXE, ASSIGN-LR-NEXE and ASSIGN-NLR into corresponding applications of ASSIGN. \square

PROPOSITION C.11. *Let \mathcal{P} be a set of programs, e be an expression and v be a variable, such that for any $p_1, p_2 \in \mathcal{P}$, there is program p with locations l_1, l_2 such that $p_1 = p[l_1/v := e; l_1]$ and $p_2 = p[l_2/v := e; l_2]$, $t := (\sigma_{in}, \phi)$ be a test. Then, $\overset{t}{\sim}_{deps}$ as in Definition 3.7 is an equivalence relation in \mathcal{P} .*

PROOF. **Reflexivity.** Let $p \in \mathcal{P}$ be a program, s is the location of “ $v := e$ ” in p , $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$. $\langle p, \sigma_{in}, \{s\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{s\}, _, _ \rangle$ (Lemma C.8). Consequently, $p \overset{t}{\sim}_{deps} p$. **Symmetry.** Let p_1, p_2 be programs such that there is program p with locations l_1, l_2 such that $p_1 = p[l_1/v := e; l_1]$, $p_2 = p[l_2/v := e; l_2]$, and $p_1 \overset{t}{\sim}_{deps} p_2$. Since $\langle p_1, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{l_1, l_2\}, _, _ \rangle$, $\langle p_2, \sigma_{in}, \{l_1, l_2\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{l_1, l_2\}, _, _ \rangle$ (Lemma C.9). Therefore $p_2 \overset{t}{\sim}_{deps} p_1$. **Transitivity.** Let p_1, p_2, p_3 be programs such that there is program p with locations l_1, l_2, l_3 such that $p_1 = p[l_1/v := e; l_1]$, $p_2 = p[l_2/v := e; l_2]$ and $p_3 = p[l_3/v := e; l_3]$, and $p_1 \overset{t}{\sim}_{deps} p_2$, $p_1 \overset{t}{\sim}_{deps} p_3$. Therefore, $\langle p_1, \sigma_{in}, \{l_1, l_2, l_3\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{l_1, l_2, l_3\}, _, _ \rangle$. Therefore, the locations $\{l_1, l_2, l_3\}$ belongs to each test-equivalent interval of the corresponding trace. Consequently, $\langle p_2, \sigma_{in}, \{l_2, l_3\}, \emptyset, false \rangle \Downarrow_{deps} \langle _, \{l_2, l_3\}, _, _ \rangle$. Therefore, $p_2 \overset{t}{\sim}_{deps} p_3$ \square

PROPOSITION C.12. *Let p_1 and p_2 be programs such that there is program p with locations l_1, l_2 such that $p_1 = p[l_1/v := e; l_1]$ and $p_2 = p[l_2/v := e; l_2]$, $t := (\sigma_{in}, \phi)$ be a test. If $p_1 \overset{t}{\sim}_{deps} p_2$, then p_1 and p_2 either both pass t or both fail t .*

PROOF. Without loss of generality we assume that the last statement of p_1 and p_2 is $v := v$. Then, p_1 and p_2 produce the same output states for the inputs σ_{in} (Lemma C.9 and Lemma C.10). Therefore, p_1 and p_2 either both pass t or both fail t . \square

PROPOSITION C.13. *The relation $\overset{t}{\sim}_{deps}$ is a test-equivalence relation according to Definition 1.1.*

PROOF. Proposition C.11 states that $\overset{t}{\sim}_{deps}$ is an equivalence relation and Proposition C.12 states that any two programs equivalent w.r.t. $\overset{t}{\sim}_{deps}$ produce indistinguishable results for the given test. Therefore, $\overset{t}{\sim}_{deps}$ is a test-equivalence relation. \square

C.3 Composition of test-equivalence relations

This section describes properties and their proofs for the semantics in Figure 13. In an analogous manner to the previous sections, we define the function *Trace* for the dependency-based semantics in such a way that it transforms each rule application into a function $f_{RULE}^i : \Sigma \times (Stmt \rightarrow$

$2^{Expr} \times (Stmt \rightarrow 2^{Expr}) \times \mathbb{B} \rightarrow \Sigma \times (Stmt \rightarrow 2^{Expr}) \times (Stmt \rightarrow 2^{Expr}) \times \mathbb{B}$, where the seconds and the third elements of the tuple denote global and current test-equivalence classes, and the last boolean element indicates if the inserted assignment has been executed.

LEMMA C.14. *Let p be a program, e be an expression and v be a variable, such that the assignment “ $v := e$ ” is at the location s of p , $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$. Then, $\langle p, \sigma_{in}, l_{in}, \lambda x. \emptyset, false \rangle \Downarrow_* \langle _, l_{in}, _, _ \rangle$, where $l_{in} := \lambda x. \text{if } x = s \text{ then } \{e\} \text{ else } \emptyset$.*

PROOF. Consider a trace of the derivation tree of the relation $\langle p, \sigma_{in}, l_{in}, \lambda x. \emptyset, false \rangle \Downarrow_* \langle _, _, _, _ \rangle$:

$$\{ f_{R_0}, f_{R_1}, \dots, f_{R_k} \}$$

For each i , if $\langle _, l, c, x \rangle = f_{R_i} \circ \dots \circ f_{R_k}(\sigma_{in}, l_{in}, \lambda x. \emptyset, false)$, then $x = true \Leftrightarrow e \in c(s)$ since x is set to *true* only when the rule ASSIGN-INS is applied (that updates c for the argument s), and x is set to *false* in the rules VAR-LEFT-EXE, VAR-LEFT-NEXE, ASSIGN-LR-EXE and ASSIGN-LR-NEXE (that make $c = \lambda x. \emptyset$). Therefore, when the rules VAR-LEFT-NEXE and ASSIGN-LR-NEXE are applied, $e \notin c(s)$. Thus, $l(s)$ always includes at least e in the derivation tree. Therefore, $\langle p, \sigma_{in}, l_{in}, \lambda x. \emptyset, false \rangle \Downarrow_* \langle _, l_{in}, _, _ \rangle$. \square

LEMMA C.15. *Let p_1 and p_2 be programs, e_1 and e_2 be expressions and v be a variable, such that there is program p with locations l_1, l_2 so that $p_1 = p[l_1/v := e_1; l_1]$, $p_2 = p[l_2/v := e_2; l_2]$, $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$, $p_1 \stackrel{t}{\sim}_* p_2$. Assume that $I_0^1, I_1^1, \dots, I_k^1$ are the test-equivalent intervals in the trace of p_1 for a derivation tree $\langle p_1, \sigma_{in}, _, _ \rangle \Downarrow_* \langle _, _, _, _ \rangle$. Then, for each I_i^1 there is a corresponding test-equivalent interval I_i^2 in the trace of $\langle p_2, \sigma_{in}, _, _ \rangle \Downarrow_* \langle _, _, _, _ \rangle$ such the first rules in the intervals receive the same inputs state for the same program location and the last rules in the intervals produce the same state for the same program location in the trace.*

PROOF. Assume $I_i^1 := \{f_{R_0}^1, \dots, f_{R_k}^1\}$ and $I_i^2 := \{f_{R_0}^2, \dots, f_{R_l}^2\}$. We prove that for any state σ , if $\langle \sigma', _, _, _ \rangle = f_{R_0}^1 \circ \dots \circ f_{R_k}^1(\sigma, _, _, _)$ and the rules $f_{R_0}^1$ and $f_{R_0}^2$ are applied for the same program location, then $\langle \sigma', _, _, _ \rangle = f_{R_0}^2 \circ \dots \circ f_{R_l}^2(\sigma, _, _, _)$ and $f_{R_l}^1$ and $f_{R_l}^2$ are applied for the same program location. This fact is proved by induction on the rule index in the test-equivalent interval, excluding the applications of ASSIGN-INS. Then, this fact can be used to prove the lemma by induction on the index of the test-equivalent interval. \square

LEMMA C.16. *Let p be program, $t := (\sigma_{in}, \phi)$ be a test. If $\langle p, \sigma_{in}, _, _ \rangle \Downarrow_* \langle \sigma_{out}, _, _, _ \rangle$, then $\langle p, \sigma_{in} \rangle \Downarrow \sigma_{out}$.*

PROOF. A derivation tree for the semantics in Figure 13 can be transformed into a valid derivation tree of the standard semantics (Figure 6) by removing the arguments l, c and x , changing the applications of VAR-LEFT-EXE, VAR-LEFT-NEXE and VAR-NLEFT into corresponding applications of VAR and the applications of ASSIGN-INS, ASSIGN-LR-EXE, ASSIGN-LR-NEXE and ASSIGN-NLR into corresponding applications of ASSIGN, removing applications of STMT-L. \square

PROPOSITION C.17. *Let \mathcal{P} be a set of programs, e be and expression and v be a variable, such that for any $p_1, p_2 \in \mathcal{P}$, there is program p with locations l_1, l_2 and expressions e_1 and e_2 such that $p_1 = p[l_1/v := e_1; l_1]$ and $p_2 = p[l_2/v := e_2; l_2]$, $t := (\sigma_{in}, \phi)$ be a test. Then, $\stackrel{t}{\sim}_*$ as in Definition 3.7 is an equivalence relation in \mathcal{P} .*

PROOF. **Reflexivity.** Let $p \in \mathcal{P}$ be a program, s is the location of “ $v := e$ ” in p , $Inserted := \lambda s. s = “v := e”$, $Left := \lambda v'. v' = v$, $Right := \lambda v'. v' \in Var(e)$. $\langle p, \sigma_{in}, l_{in}, \lambda x. \emptyset, false \rangle \Downarrow_* \langle _, l_{in}, _, _ \rangle$ where $l_{in} := \lambda x. \text{if } x = s \text{ then } \{e\} \text{ else } \emptyset$ (Lemma C.14). Consequently, $p \stackrel{t}{\sim}_* p$. **Symmetry.** Let p_1, p_2 be programs and e_1 and e_2 be expressions such that there is program p with locations

l_1, l_2 such that $p_1 = p[l_1/v := e_1; l_1]$, $p_2 = p[l_2/v := e_2; l_2]$, and $p_1 \stackrel{t}{\sim}_* p_2$. Since $\langle p_1, \sigma_{in}, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, _ \mapsto \emptyset\}, \emptyset, false \rangle \Downarrow_* \langle _, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, _ \mapsto \emptyset\}, _ \rangle$, $\langle p_2, \sigma_{in}, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, _ \mapsto \emptyset\}, \emptyset, false \rangle \Downarrow_* \langle _, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, _ \mapsto \emptyset\}, _ \rangle$ (Lemma C.15). Therefore $p_2 \stackrel{t}{\sim}_* p_1$. **Transitivity.** Let p_1, p_2, p_3 be programs and e_1, e_2, e_3 be expressions such that there is program p with locations l_1, l_2, l_3 such that $p_1 = p[l_1/v := e_1; l_1]$, $p_2 = p[l_2/v := e_2; l_2]$ and $p_3 = p[l_3/v := e_3; l_3]$, and $p_1 \stackrel{t}{\sim}_* p_2$, $p_1 \stackrel{t}{\sim}_* p_3$. Therefore, $\langle p_1, \sigma_{in}, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, l_3 \mapsto \{e_3\}, _ \mapsto \emptyset\}, \emptyset, false \rangle \Downarrow_* \langle _, \{l_1 \mapsto \{e_1\}, l_2 \mapsto \{e_2\}, l_3 \mapsto \{e_3\}, _ \mapsto \emptyset\}, _ \rangle$. Therefore, the locations $\{l_1, l_2, l_3\}$ belongs to each test-equivalent interval of the corresponding trace. Consequently, $\langle p_2, \sigma_{in}, \{l_2 \mapsto \{e_2\}, l_3 \mapsto \{e_3\}, _ \mapsto \emptyset\}, \emptyset, false \rangle \Downarrow_* \langle _, \{l_2 \mapsto \{e_2\}, l_3 \mapsto \{e_3\}, _ \mapsto \emptyset\}, _ \rangle$. Therefore, $p_2 \stackrel{t}{\sim}_* p_3$ \square

PROPOSITION C.18. *Let p_1 and p_2 be programs and e_1 and e_2 be expressions such that there is program p with locations l_1, l_2 such that $p_1 = p[l_1/v := e_1; l_1]$ and $p_2 = p[l_2/v := e_2; l_2]$, $t := (\sigma_{in}, \phi)$ be a test. If $p_1 \stackrel{t}{\sim}_* p_2$, then p_1 and p_2 either both pass t or both fail t .*

PROOF. Without loss of generality we assume that the last statement of p_1 and p_2 is $v := v$. Then, p_1 and p_2 produce the same output states for the inputs σ_{in} (Lemma C.15 and Lemma C.16). Therefore, p_1 and p_2 either both pass t or both fail t . \square

PROPOSITION C.19. *The relation $\stackrel{t}{\sim}_*$ is a test-equivalence relation according to Definition 1.1.*

PROOF. Proposition C.17 states that $\stackrel{t}{\sim}_*$ is an equivalence relation and Proposition C.18 states that any two programs equivalent w.r.t. $\stackrel{t}{\sim}_*$ produce indistinguishable results for the given test. Therefore, $\stackrel{t}{\sim}_*$ is a test-equivalence relation. \square

Received October 2017; revised May 2018; accepted July 2018