

# Information and Software Technology

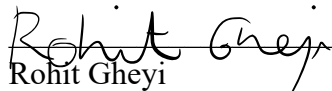
## Identifying Method-Level Mutation Subsumption Relations using Z3

--Manuscript Draft--

<b>Manuscript Number:</b>	INFSOF-D-20-00199R2
<b>Article Type:</b>	Research paper
<b>Keywords:</b>	Mutation Analysis, Redundant Mutants, Theorem Proving
<b>Corresponding Author:</b>	Rohit Gheyi Federal University of Campina Grande Campina Grande, BRAZIL
<b>First Author:</b>	Rohit Gheyi
<b>Order of Authors:</b>	Rohit Gheyi Marcio Ribeiro Beatriz Sousa Marcio Guimarães Leonardo Fernandes Marcelo d'Amorim Vander Alves Leopoldo Teixeira Balduino Fonseca
<b>Abstract:</b>	<p><b>Context:</b> Mutation analysis is a popular but costly approach to assess the quality of test suites. One recent promising direction in reducing costs of mutation analysis is to identify redundant mutations, i.e., mutations that are subsumed by some other mutations. A previous approach found redundant mutants manually through truth tables but it cannot be applied to all mutations. Another work derives them using automatic test suite generators but it is a time consuming task to generate mutants and tests, and to execute tests. <b>Objective:</b> This article proposes an approach to discover redundant mutants by proving subsumption relations among method-level mutation operators using weak mutation testing. <b>Method:</b> We conceive and encode a theory of subsumption relations in the Z3 theorem prover for 37 mutation targets (mutations of an expression or statement). <b>Results:</b> We automatically identify and prove a number of subsumption relations using Z3, and reduce the number of mutations in a number of mutation targets. To evaluate our approach, we modified MuJava to include the results of 24 mutation targets and evaluate our approach in 125 classes of 5 large open source popular projects used in prior work. Our approach correctly discards mutations in 75.93% of the cases, and reduces the number of mutations by 71.38%.</p> <p><b>Conclusions:</b> Our approach offers a good balance between the effort required to derive subsumption relations and the effectiveness for the targets considered in our evaluation in the context of strong mutation testing.</p>

Dear Günther Ruhe,

We are grateful for the valuable feedback and comments on the earlier version of our article (INFOSOF-D-20-00199), and we revised it accordingly. Please find below our responses to the individual issues that reviewers raised. We also include our article highlighting the changes we made.

  
 Rohit Gheyi

Department of Computing and Systems  
 Federal University of Campina Grande  
 882, Aprígio Veloso, Bodocongó  
 Campina Grande, PB, 58429-900, Brazil  
 phone: +55 83 2101-1122, extension 2202  
 e-mail: [rohit@dsc.ufcg.edu.br](mailto:rohit@dsc.ufcg.edu.br)

Reviewer comments

[Our answers](#)

### Reviewer #1:

Next we define the kills relation in Definition 1. -> We now define the kills relation  
[Fixed \(see Section 2\).](#)

The scope of quantification over  $m_1, m_2$  is odd in Definition 2: "Consider a program  $p$  and two distinct mutations,  $M_1$  and  $M_2$ , that are applied on the same mutation target. For all targets  $tgt$  in  $p$  and all mutants  $m_1$  and  $m_2$  generated from  $M_1$  and  $M_2$ , respectively, on  $tgt$ , we say that  $M_2$  subsumes  $M_1$  iff: ...."

The intent seems to be that the definition is about " $M_2$  subsumes  $M_1$ ". Then the auxiliary entities,  $m_1, m_2, tgt$  are used to define this concept, they are not part of what is defined. So on the left side of the definition you would have " $M_2$  subsumes  $M_1$ ", on the right side, all the auxiliary entities:

"Consider a program  $p$  and two distinct mutations,  $M_1$  and  $M_2$ , that are applied on the same mutation target. We say that  $M_2$  subsumes  $M_1$  iff for all targets  $tgt$  in  $p$  and all mutants  $m_1$  and  $m_2$  generated from  $M_1$  and  $M_2$ , respectively, on  $tgt$ : ...."

[Fixed \(see Section 2\).](#)

I found this snippet needlessly complicated:

```
def keepNonRedundantMutants(muts):
    redundant = set ()
    for i in range (len(muts)) :
        for j in range (len(muts)) :
            if i != j and redundantMutants (muts[i], muts[j]) :
                if muts [i] not in redundant and muts [j] not in redundant :
```

```

    redundant.add(muts [j])
return [m for m in muts if m not in redundant ]
It would be equivalent to:
def keepNonRedundantMutants(muts):
    non_redundant = []
    for m1 in muts:
        if not any(redundantMutants(m1, m2) for m2 in non_redundant):
            non_redundant.append(m1)
    return non_redundant
We updated Listing 7.

```

### Reviewer #3:

I am still unclear about the main issue with subsumption: what is its relation to semantic equivalence? again, the most natural way to minimize redundancy is to select a single representative from each equivalence class (modulo semantic equivalence) of the set of mutants. the authors say that removing subsumed mutants is the same thing as selecting a rep from each equivalence class, but that is patently false, if only for the following reasons: first equivalence is an equivalence relation whereas subsumption is an ordering relation; second, equivalence is an intrinsic property of the mutants whereas subsumption involves the two mutants and a mutation target. At the very least, I would change the second condition of definition to a  $\backslash\text{subsepeq}$  relation rather than a  $\backslash\text{subset}$  relation, and analyze the relation between semantic equivalence and mutual subsumption. I believe the latter logically implies the former, which means subsumption gives a superset of the set of mutants provided by equivalence.

We followed your suggestion and changed the general definition to a  $\backslash\text{subsepeq}$  relation rather than a  $\backslash\text{subset}$  relation. We changed Definition 2, Listing 4, the last paragraph of Section 2, and the fourth paragraph of Section 4.1. This modification does not impact on our results presented in Table 1. The only expected change is that now some subsumption relation graphs may have bidirectional arrows between two distinct  $m_1$  and  $m_2$  nodes when  $\text{kills}(p, m_1) = \text{kills}(p, m_2)$ .

also the authors did not answer my question whether minimal means minimal for inclusion or having a minimal cardinality (I believe it is the former, hence you may end up with a larger set than needed).

It is having a minimal cardinality. Consider we have three mutations  $m_1$ ,  $m_2$ , and  $m_3$ , such that:

- $\text{kills}(p, m_1) \neq \emptyset$ ;
- $\text{kills}(p, m_1) \subset \text{kills}(p, m_3)$ ;
- $\text{kills}(p, m_1) = \text{kills}(p, m_2)$ .

Developers can use the following minimal sets:  $\{m_1\}$  or  $\{m_2\}$ . We include a clarification in the end of the second paragraph of Section 4.

Definition 1 is incomplete because it fails to specify what happens when one of the programs or both programs fail to terminate.

We updated Definition 1 to state that we assume  $p$  and  $m$  always terminate when running any test case.

also I asked whether subsumption is binary or higher order because it is conceivable that, e.g.  $M1$  does not subsume  $M2$ ,  $M3$  does not subsume  $M2$ , but the combination of  $M1$  and  $M3$  subsume  $M2$ . in which case you want to include  $M1$  and  $M3$  and exclude  $M2$ . is that possible? are you making provisions? do they all have to have the same mutation target?

We encoded some higher order mutations proposed in the literature for the same mutation target in our theory (see Listing 8). In our preliminary results, we do not have examples that fit this scenario:

- $M1$  does not subsume  $M2$ ,  $M3$  does not subsume  $M2$ , but the combination of  $M1$  and  $M3$  subsume  $M2$ .

As future work, we intend to investigate more about higher order mutants and see whether this scenario occurs.

## Identifying Method-Level Mutation Subsumption Relations using Z3

**Context:** Mutation analysis is a popular but costly approach to assess the quality of test suites. One recent promising direction in reducing costs of mutation analysis is to identify redundant mutations, i.e., mutations that are subsumed by some other mutations. A previous approach found redundant mutants manually through truth tables but it cannot be applied to all mutations. Another work derives them using automatic test suite generators but it is a time consuming task to generate mutants and tests, and to execute tests.

**Objective:** This article proposes an approach to discover redundant mutants by proving subsumption relations among method-level mutation operators using weak mutation testing.

**Method:** We conceive and encode a theory of subsumption relations in the Z3 theorem prover for 37 mutation targets (mutations of an expression or statement).

**Results:** We automatically identify and prove a number of subsumption relations using Z3, and reduce the number of mutations in a number of mutation targets. To evaluate our approach, we modified MuJava to include the results of 24 mutation targets and evaluate our approach in 125 classes of 5 large open source popular projects used in prior work. Our approach correctly discards mutations in 75.93% of the cases, and reduces the number of mutations by 71.38%.

**Conclusions:** Our approach offers a good balance between the effort required to derive subsumption relations and the effectiveness for the targets considered in our evaluation in the context of strong mutation testing.

# Identifying Method-Level Mutation Subsumption Relations using Z3<sup>\*</sup>

Rohit Gheyi<sup>a,\*</sup>, Márcio Ribeiro<sup>b</sup>, Beatriz Sousa<sup>a</sup>, Marcio Guimarães<sup>b</sup>, Leo Fernandes<sup>c</sup>, Marcelo d'Amorim<sup>d</sup>, Vander Alves<sup>e</sup>, Leopoldo Teixeira<sup>d</sup>, Balduino Fonseca<sup>b</sup>

<sup>a</sup>*Department of Computing and Systems, UFCG, Campina Grande-PB, Brazil*

<sup>b</sup>*Computing Institute, UFAL, Maceió-AL, Brazil*

<sup>c</sup>*IFAL, Maceió-AL, Brazil*

<sup>d</sup>*Informatics Center, Universidade Federal de Pernambuco, Recife-PE, Brazil*

<sup>e</sup>*Computer Science Department, UnB, Brasília-DF, Brazil*

---

## Abstract

**Context:** Mutation analysis is a popular but costly approach to assess the quality of test suites. One recent promising direction in reducing costs of mutation analysis is to identify redundant mutations, i.e., mutations that are subsumed by some other mutations. A previous approach found redundant mutants manually through truth tables but it cannot be applied to all mutations. Another work derives them using automatic test suite generators but it is a time consuming task to generate mutants and tests, and to execute tests.

**Objective:** This article proposes an approach to discover redundant mutants by proving subsumption relations among method-level mutation operators using weak mutation testing.

**Method:** We conceive and encode a theory of subsumption relations in the Z3 theorem prover for 37 mutation targets (mutations of an expression or statement).

**Results:** We automatically identify and prove a number of subsumption rela-

---

<sup>\*</sup>This work was partially supported by CNPq and CAPES grants.

<sup>\*</sup>Corresponding Author

*Email addresses:* rohit@dsc.ufcg.edu.br (Rohit Gheyi), marcio@ic.ufal.br (Márcio Ribeiro), beatriz.souza@ccc.ufcg.edu.br (Beatriz Sousa), masg@ic.ufal.br (Marcio Guimarães), leonardo.oliveira@ifal.edu.br (Leo Fernandes), damorim@cin.ufpe.br (Marcelo d'Amorim), valves@unb.br (Vander Alves), lmt@cin.ufpe.br (Leopoldo Teixeira), balduino@ic.ufal.br (Balduino Fonseca)

tions using Z3, and reduce the number of mutations in a number of mutation targets. To evaluate our approach, we modified MUJAVA to include the results of 24 mutation targets and evaluate our approach in 125 classes of 5 large open source popular projects used in prior work. Our approach correctly discards mutations in 75.93% of the cases, and reduces the number of mutations by 71.38%.

**Conclusions:** Our approach offers a good balance between the effort required to derive subsumption relations and the effectiveness for the targets considered in our evaluation in the context of strong mutation testing.

*Keywords:* Mutation Analysis, Redundant Mutants, Theorem Proving

---

## 1. Introduction

Mutation analysis is a popular technique to assess quality of test suites [6, 40, 46]. The technique introduces variations in code and checks if those variations are observable through test execution. Applying a mutation to a program yields a mutant. A mutant is said to be killed if a test case in the test suite fails on a given mutant; a mutant is said to survive otherwise. The intuition is that a test suite that kills more mutants is more adequate to detect defects when they actually occur [20].

Usually, the costs of using mutation analysis are high, mainly due to the high number of generated mutants and the high computing time to execute the test suite against each mutant. However, some mutants are redundant, that is, they may not be necessary for the effectiveness of mutation analysis and thus we may discard them [45]. We can speed up execution time using multi-execution, parallel execution, and so on. But reducing cost is still important. Redundant mutants do not contribute to the test assessment process because they are killed when other mutants are also killed [26, 45]. Redundant mutants are always subsumed by other mutants. The generation of these mutants increases the total cost and does not help to improve effectiveness of the test suite. Ammann et al. [1] empirically identified that a number of the generated mutants are redundant. Also, Papadakis et al. [47] identified that such redundant mutants

inflate the mutation score and that a number of recent research papers are vulnerable to threats to validity due to the effect of these mutants.

To identify redundant mutants, we can take subsumption relations into account. Kaminski et al. [24] manually constructed subsumption hierarchies with the support of truth tables produced by the outcomes of mutants associated with the *Relational Operator Replacement* (ROR) mutation operator. This operator generates seven different mutations, but Kaminski et al. [24] identified that only three mutations are sufficient to cover all input domains, yielding a reduction of 57% of redundant mutants. Just et al. [21] expanded this idea with two more mutation operators. Both works use truth tables to infer logical relationships across the operations. Although the idea is promising, we cannot apply it for non-logical operators. For instance, a binary expression with two numeric variables `a + b` has a very large set of input possibilities, which turns the manual and logical approach more difficult. Guimarães et al. [11] proposed an approach to yield dynamic subsumption relations among method-level mutants by using automatic test suite generators, such as Randoop [44] and EvoSuite [7] in the context of strong mutation testing. However, the approach is time consuming since it needs to generate mutants, compile them, generate test suites, and execute them.

In this article, we propose an approach consisting of six steps to discover subsumption relations among method-level mutations using theorem proving in the context of weak mutation testing [13]. We encode a theory of subsumption relations in Z3 and use its theorem prover [36] to automatically identify redundant mutations (Section 4). We consider most of the method-level mutation operators available in the MUJAVA tool [32, 33]. We reduce the number of mutations in a number of mutation targets (mutations of an expression or statement). A mutation target is a language expression or statement in which it is possible to apply a set of mutations of one or more mutation operators [11].

To evaluate our approach, we modify MUJAVA to include the results of 24 mutation targets and evaluate our approach in 125 classes of 5 real projects. Our approach achieves an effectiveness (the percentage of mutants correctly



discarded by our technique) of 75.93% and a reduction rate (the percentage of mutants discarded by our technique) of 71.38%. We achieve a good cost-benefit ratio between the effort required to derive the mutation subsumption relations and the effectiveness for the targets considered in our evaluation in the context of strong mutation testing. Moreover, we show that the random sampling strategy requires a sampling rate greater than 60% to achieve a similar effectiveness of our approach. So, our reduction mutation strategy is not considered harmful.

We organize this article as follows. We explain mutant subsumption relations in Section 2, and present a motivating example in Section 3. Section 4 describes our approach to identify subsumption relations using Z3. Section 5 presents the evaluation of our approach. Finally, we relate our approach to others (Section 6), and present concluding remarks (Section 7).

## 2. Mutation Subsumption Relations

Mutation analysis uses mutation operators to introduce faults in the program to create mutants deliberately [6]. In this context, there is a wide variety of mutation operators. Each mutation operator can implement a set of mutations. In this work, we follow the same definition for “mutation” of previous work [23]: a *mutation* refers to a syntactic change (e.g.,  $a \ \&\& \ b \mapsto a \ || \ b$ ).

Subsumption relations identify redundancy in sets of mutations and hence can be used to optimize approaches to both mutant and test generation [27]. The subsumed mutants do not need to be generated, and test generation methods can target subsuming mutants.

We now define the **kills** relation.

**Definition 1.** *Consider a program  $p$ . We apply a mutation  $M$  to  $p$  and yield one or more mutants. Let  $m$  be one of them. Both  $p$  and  $m$  always terminate when running any test case. We define the  $\text{kills}(p, m)$  function that yields all test cases that have different return values in  $p$  and  $m$ .*

For example, consider the  $p = x+y$  program. Suppose we apply a mutation  $M$  converting the arithmetic operator  $+$  to the arithmetic operator  $-$ . We yield the

$m = x-y$  mutant. In this example,  $kills(p, m)$  yields a non empty set of test cases. A test case assigns values to all variables in  $p$ . It contains the following test case  $t=(x=1, y=1)$  that yields different values in  $p$  (2) and  $m$  (0).

We define the subsumption relation in Definition 2.

**Definition 2.** *Consider a program  $p$  and two distinct mutations,  $M_1$  and  $M_2$ , that are applied on the same mutation target. We say that  $M_2$  subsumes  $M_1$  iff for all targets  $tgt$  in  $p$  and all mutants  $m_1$  and  $m_2$  generated from  $M_1$  and  $M_2$ , respectively, on  $tgt$ :*

1.  $kills(p, m_2) \neq \emptyset$
2.  $kills(p, m_2) \subseteq kills(p, m_1)$

The first condition of Definition 2 guarantees that  $m_2$  is not an equivalent mutant [34]. The program and the mutant have at least one test case that yields different values. In the second condition of Definition 2, the set of test cases that kills  $m_2$  is a subset of the set of test cases that kill  $m_1$ . Notice that we can have more test cases that kill  $m_1$  but cannot kill  $m_2$ . In this way, it is easier to kill  $m_1$  than  $m_2$ . So, we say that  $m_2$  subsumes  $m_1$ . We do not need to generate  $m_1$  during mutation testing. Studying mutation subsumption relation can help us build more efficient mutation testing tools, significantly improving the applicability of mutation testing in industry by helping to minimize one of the challenges [49].

### 3. Motivating Example

Consider a binary expression with a relational operator  $lexp <op> rexp$ , where  $lexp$  and  $rexp$  indicate expressions or literals and  $<op>$  is a relational operator ( $=$ ,  $!=$ ,  $>$ ,  $>=$ ,  $<$ , or  $<=$ ). The Relational Operator Replacement (ROR) mutation operator performs seven mutations, replacing the original operator  $<op>$  with each of the other five relational operators and replacing the entire expression with **true** and **false**. Thus, for the binary expression  $a > b$ , the ROR operator performs the following seven mutations [11]:

1.  $a > b \mapsto a == b;$
2.  $a > b \mapsto a != b;$
3.  $a > b \mapsto a >= b;$
4.  $a > b \mapsto a < b;$
5.  $a > b \mapsto a <= b;$
6.  $a > b \mapsto \text{true};$
7.  $a > b \mapsto \text{false}.$

However, some mutations may not be necessary for the effectiveness of mutation analysis and are actually useless. An equivalent mutant is syntactically different from the original program but has the same semantics [34]. In this work, we focus on redundant mutants. To identify them, we rely on subsumption relations, as defined in Section 2.

For instance, consider the binary expression  $a > b$  and two mutants:  $a >= b$  and  $a <= b$ . Notice that both mutants are not equivalent to the original binary expression using weak mutation testing. If  $a$  is different from  $b$  in a test case, we kill  $a <= b$  but we cannot kill  $a >= b$ . If  $a$  is equal to  $b$  in a test case, we kill both mutants. Since (i) all test cases that kill  $a >= b$  also kill  $a <= b$ , and (ii) there are some test cases that kill  $a <= b$  but cannot kill  $a >= b$ , we conclude that ROR ( $>=$ ) subsumes ROR ( $<=$ ) for the mutation target  $a > b$ . As a consequence, we must not apply ROR ( $<=$ ) in this mutation target if we apply ROR ( $>=$ ) using weak mutation testing, hence reducing the number of redundant mutants.

Previous works manually found redundant mutants through the truth table [24, 21]. Although the idea is promising, it can only be applied for logical and relational operators. Guimarães et al. [11] used automatic test generation to identify subsumption relations using strong mutation testing. However, we cannot have full confidence in the results derived using automatic test suite generators for some types of variables, such as integer numbers. Moreover, it is a time consuming approach to derive the subsumption relations using automatic test suite generators. In this work, we focus on proposing an approach

to automatically derive sound method-level mutation subsumption relations in a theorem prover using weak mutation testing.

#### **4. Encoding and Proving Subsumption Relations**

In this section, we propose a technique to prove subsumption relations using weak mutation testing. We focus on code fragments. We use the Z3 [36] API for Python, which has a theorem prover. We consider most MUJAVA method-level mutation operators [33], such as operators that mutate arithmetic, relational, and logical expressions, and variable assignment statements. We do not focus on the object-oriented ones, i.e., the class-level mutation operators.

Table 1: It presents the mutation targets, method-level mutations that each operator is able to create in the corresponding target, a minimal set of mutations for each target identified in our approach, and the size of a minimal set of mutations compared to the original one. `OP1`: select CDL, ODL, or VDL. We use the following variables. `exp`: unary expression, such as identifiers, variables, literals; `lexp` and `rexp`: unary expressions, or binary expression; `lhs`: identifiers, or variables used in statements; `rhs`: unary expressions, or binary expression used in statements.

Mutation Target	Mutation Operators	Minimal Set of Mutations	Size
<code>lexp + rexp (for Z<sup>+</sup>)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	AORB(*)	12.5%
<code>lexp + rexp (for Z)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	All	100%
<code>lexp - rexp (for Z<sup>+</sup>)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp)	12.5%
<code>lexp - rexp (for Z)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	All	100%
<code>lexp * rexp (for Z<sup>+</sup>)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	AORB(+), OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp)	87.5%
<code>lexp * rexp (for Z)</code>	AORB (2), VDL (2), CDL (2), ODL (2)	All	100%
<code>lexp ^ rexp (bool)</code>	COR (4), ROR(2), COI (3), VDL (2), CDL (2), ODL (2)	COR(False), COR(!)	13.3%
<code>lexp &amp;&amp; rexp</code>	COR (4), ROR(2), COI (3), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp), ROR(==), COR(False)	26.7%
<code>lexp    rexp</code>	COR (4), ROR(2), COI (3), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp), ROR(!=), COR(True)	26.7%
<code>lexp == rexp (bool)</code>	ROR (1), COI (3), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp)	20%
<code>lexp != rexp (bool)</code>	ROR (1), COI (3), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp)	20%
<code>lexp == rexp</code>	ROR (7), COI (1)	ROR(False), ROR(>=), ROR(<=)	37.5%
<code>lexp != rexp</code>	ROR (7), COI (1)	ROR(<), ROR(True), ROR(>)	37.5%
<code>lexp &gt; rexp</code>	ROR (7), COI (1)	ROR(False), ROR(!=), ROR(>=)	37.5%
<code>lexp &gt;= rexp</code>	ROR (7), COI (1)	ROR(True), ROR(==), ROR(>)	37.5%
<code>lexp &lt; rexp</code>	ROR (7), COI (1)	ROR(False), ROR(!=), ROR(<=)	37.5%
<code>lexp &lt;= rexp</code>	ROR (7), COI (1)	ROR(True), ROR(==), ROR(<)	37.5%
<code>lexp != rexp (obj)</code>	ROR (7), COI (1)	ROR(True), ROR(>), ROR(<)	37.5%
<code>lexp &amp; rexp</code>	LOR (2), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp)	25%
<code>lexp   rexp</code>	LOR (2), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp), LOR( ^ )	37.5%
<code>lexp ^ rexp</code>	LOR (2), SOR (2), CDL (2), ODL (2)	LOR(!)	12.5%
<code>lexp &gt;&gt; rexp</code>	LOR (3), SOR (1), VDL (2), CDL (2), ODL (2)	OP <sub>1</sub> (lexp), OP <sub>1</sub> (rexp), LOR( ^ ), LOR(!), LOR(&), SOR(<<)	60%
<code>lexp &lt;&lt; rexp</code>	LOR (3), SOR (1), VDL (2), CDL (2), ODL (2)	LOR( ^ ), LOR(&), SOR(>>)	30%
<code>exp</code>	AOIS (4), AOIU (1), LOI (1)	AOIU(-exp)	16.7%
<code>+exp</code>	AODU (1), LOI (1), ODL (1)	LOI(-exp)	33.3%
<code>-exp</code>	AODU (1), LOI (1), ODL (1)	AODU(exp)	33.3%
<code>++exp</code>	AORS (1), AODS (1), LOI (1), ODL (1)	AODS(exp), LOI(-exp)	50%
<code>exp++</code>	AORS (1), AODS (1), LOI (1), ODL (1)	LOI(-exp)	25%
<code>--exp</code>	AORS (1), AODS (1), LOI (1), ODL (1)	AODS(exp), LOI(-exp)	50%
<code>exp--</code>	AORS (1), AODS (1), LOI (1), ODL (1)	LOI(-exp)	25%
<code>!exp</code>	COD (1), ODL (1)	COD(exp)	50%
<code>~exp</code>	AODU (1), LOD (1), ODL (1)	LOD(exp)	33.3%
<code>lhs += rhs (for Z<sup>+</sup>)</code>	ASRS (2), ODL (1), SDL (1)	ASRS(*=)	25%
<code>lhs -= rhs (for Z<sup>+</sup>)</code>	ASRS (2), ODL (1), SDL (1)	ODL(lhs=rhs)	25%
<code>lhs *= rhs (for Z<sup>+</sup>)</code>	ASRS (2), ODL (1), SDL (1)	ODL(lhs=rhs), ASRS(+), SDL	75%
<code>lhs &lt;&lt;= rhs</code>	ASRS (1), ODL (1), SDL (1)	ASRS(>>=)	33.3%
<code>lhs &gt;&gt;= rhs</code>	ASRS (1), ODL (1), SDL (1)	All	100%
<code>lhs &amp;= rhs</code>	ASRS (2), ODL (1), SDL (1)	ODL(lhs=rhs), SDL	50%
<code>lhs  = rhs</code>	ASRS (2), ODL (1), SDL (1)	ODL(lhs=rhs), ASRS(^=), SDL	75%
<code>lhs ^= rhs</code>	ASRS (2), ODL (1), SDL (1)	ASRS(!=)	33.3%

Table 1 illustrates a number of method-level mutation targets (code fragments) in which MUJAVA is able to apply a set of mutations from one or more mutation operators. Consider the first column of the table. The first row of

the table focuses on the mutation target `lexp + rexp (for Z+)`, where `lexp` and `rexp` are positive integer expressions. In the second row of Table 1, `lexp` and `rexp` are integer expressions. For other mutation targets, we also consider boolean expressions or objects in other targets. In the second column of the table, we present the mutation operators that can be applied to each mutation target. We can apply four mutation operators in MUJAVA to the mutation target presented in the first row of the table: AORB, VDL, CDL, and ODL. Table 2 describes the mutation operators considered in our work [31]. The mutation operators can generate eight mutants, two for each operation. We provide the number of possible mutations (in parentheses) that such operator can apply into the target. So, we have in the second column of Table 1: AORB (2), VDL (2), CDL (2), ODL (2). By using our technique explained in Section 4.2, we yield a minimal set of mutations presented in the third column of Table 1. For the mutation target `lexp + rexp (for Z+)`, we only have one mutation: AORB using the operator `*`. The other seven mutants presented in the second column of Table 1 are redundant or equivalent. Since the redundant mutants do not contribute to the test assessment process because they are killed when other mutants are also killed [26, 45], our technique detects and removes them. So, we define a minimal set of mutations. [We may have more than one minimal sets, but all of them have the same set cardinality. A developer can select any of them.](#) Finally, the last column of Table 1 indicates the size of a minimal set of mutations compared to the original one. Since our minimal set for this mutation target contains one out of eight mutations, the size is 12.5%.

Table 2: Description of mutation operators.

Operator	Description
AORB	Binary Arithmetic Operator Replacement
AORS	Short-Cut Arithmetic Operator Replacement
AOIU	Unary Arithmetic Operator Insertion
AOIS	Short-Cut Arithmetic Operator Insertion
AODU	Unary Arithmetic Operator Deletion
AODS	Short-Cut Arithmetic Operator Deletion
ROR	Relational Operator Replacement
COR	Conditional Operator Replacement
COI	Conditional Operator Insertion
COD	Conditional Operator Deletion
SOR	Shift Operator Replacement
LOR	Logical Operator Replacement
LOI	Logical Operator Insertion
LOD	Logical Operator Delete
ASRS	Short-Cut Assignment Operator Replacement
SDL	Statement Deletion
VDL	Variable Deletion
CDL	Constant Deletion
ODL	Operator Deletion

This section is organized as follows. Section 4.1 presents some auxiliary functions. Section 4.2 defines the main steps of our technique. Section 4.3 encodes our technique in Z3. Finally, we present our lessons learned in Section 4.4.

#### 4.1. Auxiliary Functions

Listing 1 specifies how to prove a theorem using the Z3 Python API. It can yield three answers: the theorem is valid, invalid, or it does not know the answer. The command `Solver` creates a general purpose solver in Z3 [36]. Constraints can be added using the `add` function. The `Solver.check` method solves the constraints. The result is `sat` (satisfiable) if a solution was found. The result is `unsat` (unsatisfiable) if no solution exists. Finally, a solver may fail to solve a system of constraints and `unknown` is returned.

Listing 1: Proving a theorem in Z3.

```
def prove(theorem):  
    s = Solver()  
    s.add(Not(theorem))  
    r = s.check()  
    if r == unsat:  
        return 1 # theorem is valid  
    elif r == unknown:  
        return 2 # Z3 doesn't know the answer  
    else:  
        return 0 # theorem is invalid
```

Listing 2 presents two functions checking whether constraints are satisfiable (`isSat`) or unsatisfiable (`isUnsat`).

Listing 2: Checking constraints in Z3.

```
def check(f):  
    s = Solver()  
    s.add(f)  
    r = s.check()  
    if r == unknown:  
        print('unexpected unknown result for ', f)  
    return r  
  
def isSat(f):  
    return check(f) == sat  
  
def isUnsat(f):  
    return check(f) == unsat
```

Before specifying the subsumption relation, we encode the `kills` function in Listing 3. It defines a formula stating that `p` and `m` have different values.



Listing 3: The function kills.

```
def kills(p, m):
    return p != m
```

The `subsumption` function presented in Listing 4 checks whether a mutation subsumes another one, when considering the input program `p`. We may add some conditions (`conds`) when checking a theorem, such as restricting that all integer numbers are positive. The first part of Definition 2 states that  $\text{kills}(\mathbf{p}, \mathbf{m}_2) \neq \emptyset$ . To encode it in Z3, we define the `isNonEmpty` function, which tries to find a test case for `kills(p, m)`. To check the second condition presented in Section 2, we define the `isSubset` function, which checks whether there is no test case that is valid for `kills(p, m2)` but it is not valid for `kills(p, m1)`.

Listing 4: Defining a theorem in Z3.

```
def isNonEmpty(p, cond, m):
    return isSat(And(cond, kills(p, m)))

def isSubset(p, cond, m1, m2):
    return isUnsat(And(cond, kills(p, m2), Not(kills(p, m1))))

def subsumption(p, m1, m2, conds):
    t1 = isNonEmpty(p, conds, m2)
    t2 = isSubset(p, conds, m1, m2)
    if t1 == 1 and t2 == 1:
        return (m2, m1) # m2 subsumes m1
    else:
        return None
```

To make it easier to compare all mutations, we define the `identifySubsumptions` function (see Listing 5) that compares all possible combinations to identify whether a mutation subsumes another one. `mut`s represents a list of mutants.

Listing 5: Identifying all subsumption relations in Z3.

```

def identifySubsumptions(p, muts, conds):
    result = []
    for i in range(len(muts)):
        for j in range(len(muts)):
            if i != j:
                s = subsumption(p, muts[i], muts[j], conds)
                if (s is not None):
                    result.append(s)
    return result

```

Moreover, we also declare the `keepNonEquivalentMutants` function that keeps only non-equivalent mutants (see Listing 6). In this way, we discard equivalent mutants from our analysis, hence satisfying the first condition of Definition 2.

Listing 6: Identifying equivalent mutants in Z3.

```

def keepNonEquivalentMutants(p, muts):
    return [m for m in muts if prove(p==m)!=1]

```

Finally, we define the `keepNonRedundantMutants` function that keeps only non redundant mutants (see Listing 7). A mutant is duplicated to another mutant when both of them have the same semantics. In this way, we discard redundant mutants.

Listing 7: [Detecting redundant mutants in Z3](#).

```
def redundantMutants(m1,m2):  
    if prove(m1==m2) == 1:  
        return True  
    else:  
        return False  
  
def keepNonRedundantMutants(muts):  
    non_redundant = []  
    for m1 in muts:  
        if not any(redundantMutants(m1, m2) for m2 in non_redundant):  
            non_redundant.append(m1)  
    return non_redundant
```

#### 4.2. Steps

For each mutation target, the main steps of our approach are the following:

1. Declare variables and conditions;
2. Specify a program;
3. Specify a list of mutants;
4. Identify and remove equivalent mutants;
5. Identify and remove redundant mutants;
6. Identify subsumption relations.

Only Steps 4-6 do not change for all mutation targets.

#### 4.3. Encoding

Next we follow the steps presented in Section 4.2, use the auxiliary functions presented in Section 4.1, and identify some subsumption relations for some mutation targets presented in Table 1 using weak mutation testing.

#### 4.3.1. Boolean Expressions

Next we prove subsumption relations for boolean expressions. For a boolean expression `lexp && rexp` (the eighth row of Table 1), we simplify it to `x && y`. We declare `x` and `y` as boolean variables in the Z3 Python API (Step 1) as shown in Listing 8. We can declare other types of variables in Z3 [36]: `Int` (integer numbers), `Bool` (boolean variables), `BitVec` (bit-vector variables), `Real` (real numbers), and so on. In our work, we use `Bool`, `Int`, and `BitVec` with 32 bits. For the `x && y` target, we do not impose any condition (see first column of Table 1). So, we declare `conds=True` in our example.

In Step 2, we specify our program. In Z3, we have the following boolean operators: `And`, `Or`, `Not`, `Implies` (implication), `If`, and so on. In our example, we use the declared variables and specify our program in Z3: `And(x,y)` (see Listing 8).

After declaring variables and a program, we specify all mutants in Step 3. According to Table 1, the binary expression representing the input program `x and y` can derive the following mutants using ODL, VDL, CDL, COR, and COI operators:

- `True (COR true);`
- `False (COR false);`
- `x or y (COR ||);`
- `x (VDL/CDL/ODL(rexp));`
- `y (VDL/CDL/ODL(lexp));`
- `not(x and y) (COI !());`
- `x == y (COR ==);`
- `x != y (COR !=);`
- `x xor y (COR ^);`
- `not(x) xor y (COR COI (!x && y));`

- `x xor not(y) (COR COI (x && !y))`.

The last two of them are higher-order mutants derived from COR and COI mutant operators [18]. We consider them to show how to encode them in our approach. Next we manually specify them using the Z3 boolean operators (see Listing 8), but this process can be automated.

Listing 8: Identify Subsumption Relations for `lexp && rexp` target.

```
# Step 1
x = Bool('x')
y = Bool('y')
conds = True

# Step 2
p = And(x,y)

# Step 3
muts = [True, False, Or(x,y), x, y,
        Not(p), equals(x,y),
        Not(equals(x,y)), xor(x,y),
        xor(Not(x),y), xor(x,Not(y))]

# Step 4
muts = keepNonEquivalentMutants(p,muts)

# Step 5
muts = keepNonRedundantMutants(muts)

# Step 6
subsumptions = identifySubsumptions(p, muts, conds)
```

Next we identify non-equivalent mutants in some targets using the `keepNonEquivalentMutants` function (Step 4). For instance, consider the `exp` mutation target. Some mutants (`exp++`, and `exp--`) are equivalent to the program `exp` in our encoding using weak mutation testing.

We can further reduce the number of mutations by checking whether there are some mutants that are redundant to other ones in Step 5. We can check this by calling the `keepNonRedundantMutants` function passing the set of mutants

yielded in Step 4. For the `lexp && rexp`, all four dominant nodes are not redundant. We find some redundant mutants for other targets, such as `--exp` target. Consider the following set of mutations: `AODS(exp)`, `AORS(exp++)`, and `ODL(exp)`. The three mutants are redundant. Since they are redundant, for the `--exp` target, we can select one of them (`AODS(exp)`, `AORS(exp++)`, and `ODL(exp)`), instead of selecting all of them.

Finally, to identify all subsumption relations in Step 6, we have to call the `identifySubsumptions` function passing `p`, `mutts`, and `conds` as parameters. Based on the output, our script automatically derives the following subsumption graph presented in Figure 1 for the mutation target `lexp && rexp`. We create a node for each mutation, and an arrow between two nodes, when a mutation subsumes another one. For example, since `COR ||` subsumes `COR true`, we specify this subsumption relation by including an arrow between the nodes. For the `lexp && rexp` mutation target, our results indicate that we only need to use the following ones: `ODL lexp`, `VDL lexp`, `CDL lexp`, `ODL rexp`, `VDL rexp`, `CDL rexp`, `COR ==`, and `COR false`. These nodes dominate the others since they do not have incoming arrows. It is important to mention that `ODL exp`, and `VDL exp` or `CDL exp` yield syntactic equivalent mutants when we are dealing with variables or constants. We only need to select one of them. So, we only need to use four mutations for the following target `x and y`.

We also automated the process of creating the graph presented in Figure 1 by using the `graphviz` Python library. Listing 9 declares the `createSubsRelationGraph` function that receives the subsumptions relations identified by `identifySubsumptions`, a list of mutations and a dictionary specifying the names for each mutation.

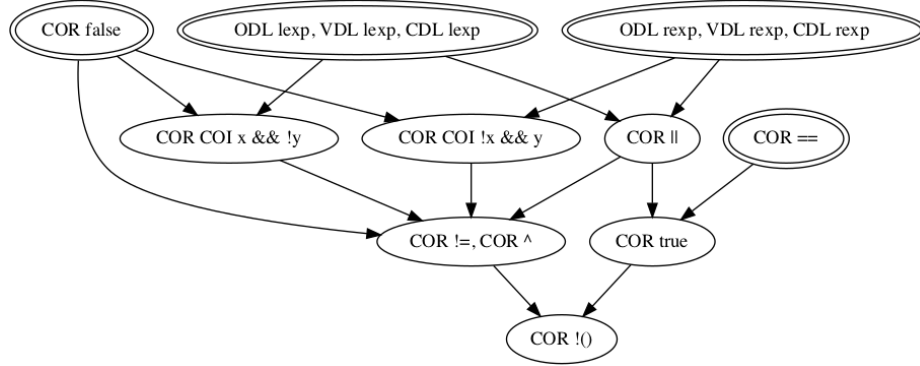


Figure 1: Mutation subsumption graph for the `lexp && rexp` mutation target. Mutations `CDL/VDL/ODL(lexp)`, `CDL/VDL/ODL(rexp)`, `COR ==`, and `COR false` dominate the other mutations.

Listing 9: Creating the subsumption relation graph.

```
def createSubRelationGraph(subsumptions, muts, mutNames):
    graph = Digraph('G')
    for m in muts:
        graph.node(mutNames[str(m)])
    for s in subsumptions:
        x = mutNames[str(s[0])]
        y = mutNames[str(s[1])]
        graph.edge(x, y)
    return graph
```

It is important to mention that `prove` (see Listing 1) does not yield `unknown` as a result in any of the results presented in Table 1. But this scenario can happen for other mutation targets. When Z3 yields `unknown`, we cannot identify subsumed relations. We recommend adding some conditions (`conds`) to the variables to avoid `unknown` in `prove`. This way, we may identify some useful subsumption relations for a restricted domain. So, we have confidence in the results given by the Z3 theorem prover. It takes a few seconds to prove all relations on a MacBook Pro 2,3 GHz Intel Core i5 with 8GB RAM memory.

#### 4.3.2. Integer Expressions

Consider the `lexp + rexp` target (the second row of Table 1). In Step 1, we declare integer variables `x` and `y` (see Listing 10). First, we will not impose any condition to identify subsumption relations (`conds=True`), since we do not have any constraint to this mutation target. We declare the `x + y` program in Python using its arithmetic operator. Then, we specify the following mutations (`AORB (2)`, `VDL (2)`, `CDL (2)`, `ODL (2)`). Since `VDL`, `CDL`, and `ODL` yield the same result, we only declare one mutation for them. In this example, we have four mutations:

- `x * y (AORB(*))`;
- `x - y (AORB(-))`;
- `x (VDL/CDL/ODL(rexp))`;
- `y (VDL/CDL/ODL(lexp))`.

Listing 10: Identify Subsumption Relations for the `lexp + rexp` target.

```
# Step 1
x = Int('x')
y = Int('y')
conds = True

# Step 2
p = x+y

# Step 3
muts = [x*y, x-y, x, y]

# Step 4
muts = keepNonEquivalentMutants(p, muts)

# Step 5
muts = keepNonRedundantMutants(muts)

# Step 6
subsumptions = identifySubsumptions(p, muts, conds)
```



We do not find equivalent and redundant mutants in Steps 4 and 5. Step 6 does not identify any subsumption relation for the `lexp + rexp` target. For instance, for the test case `x=2, y=2`, we can kill the `AORB(-)` mutation, but we cannot kill the `AORB(*)` mutation. On the other hand, for the test case `x=2, y=0`, we can kill the `AORB(*)` mutation, but we cannot kill the `AORB(-)` mutation. So, all mutations are dominant (see Figure 2) different from the result obtained by Guimarães et al. [11]. All mutation targets containing integer numbers have different subsumption graphs from Guimarães et al. [11].

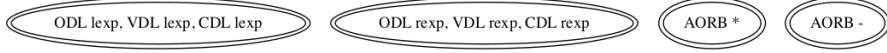


Figure 2: Mutation subsumption graph for the `lexp + rexp` mutation target. All mutations are dominant.

However, in case there are some restrictions in the developers' domain, we can reduce the number of generated mutants. For instance, suppose that all numbers are positive (for  $Z^+$ ) in the developers' domain (see first row of Table 1). In Step 1, we can specify it (`conds = And(x>0,y>0)`) using Z3 and Python operators. We can execute all steps again, and now it yields the subsumption relation presented in Figure 3. In this setting, the `AORB(*)` mutation dominates all other mutations for the `lexp + rexp` mutation target. Even considering this condition, the subsumption graph is different from the result obtained by Guimarães et al. [11], which has a limitation in their technique due to limitations in using automatic test suite generators.

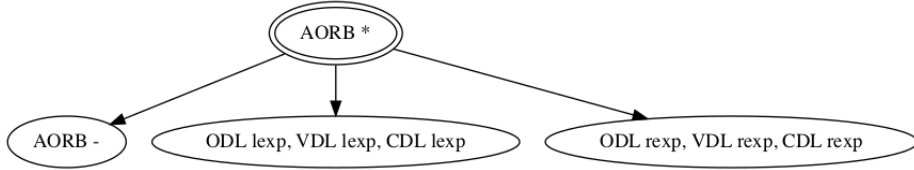


Figure 3: Mutation subsumption graph for the `lexp + rexp` mutation target considering positive integer numbers (for  $Z^+$ ).

#### *4.3.3. Expressions containing bitwise operators*

For expressions using bits, we follow a similar approach. We declare the variables `x` and `y` as a `BitVec` with 32 bits. Then we follow the same steps. It is important to mention that all bitwise operators presented in Table 1 (Mutation target column) have an equivalent bitwise operator in Python.

#### *4.3.4. Expressions containing assignment operators*

For expressions containing assignment operators, consider the `lhs &= rhs` target (see second to the last row of Table 1). The encoding is equivalent to the one presented in Section 4.3.2. In Step 1, we declare `BitVec` variables `x` and `y` containing 32 bits (see Listing 11). We will not impose any condition to identify subsumption relations (`conds=True`), since we do not have any constraint to this mutation target. The only difference between encoding expressions and commands is that we update the variable `lhs`. Since in our program and in all mutants we only have one variable `lhs` being updated, we do not need to specify it in our encoding. In summary, we encode commands in the same way we encode expressions.

Listing 11: Identify Subsumption Relations for the `lhs &= rhs` target.

```
# Step 1
x = BitVec('x',32)
y = BitVec('y',32)
conds = True
# Step 2
p = x&y
# Step 3
muts = [y,x|y,x^y,x]
# Step 4
muts = keepNonEquivalentMutants(p,muts)
# Step 5
muts = keepNonRedundantMutants(muts)
# Step 6
subsumptions = identifySubsumptions(p, muts, conds)
```

We declare the `x &= y` program in Python using its bitwise operator. Then, we specify all mutations using the Python bitwise operators:

- `y (ODL (lhs=rhs));`
- `x|y (ASRS(|=));`
- `x^y, (ASRS(^=));`
- `x (SDL).`

For the Statement Deletion mutation operator (**SDL**), it always yields **x**. Finally, we execute Steps 4-6 and it yields the subsumption relation graph presented in Figure 4.

#### 4.4. Lessons Learned

Guimarães et al. [11] used automatic test generators to derive subsumption relations using strong mutation testing. In this work, we use theorem proving in

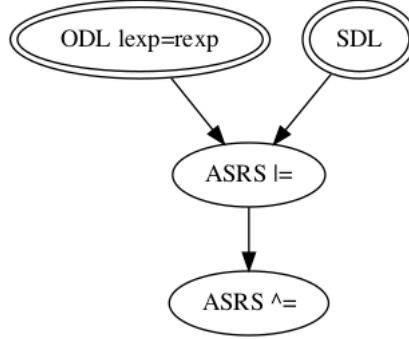


Figure 4: Mutation subsumption graph for the `lhs &= rhs` mutation target.

the context of weak mutation testing. Since all proofs are automatically done by the Z3 theorem prover, it is easier and faster to derive the subsumption relations using the approach presented here. In the approach proposed by Guimarães et al. [11], we have to generate a number of mutants using MUJAVA, compile all of them, generate a number of tests for them, and then analyze all results. It is a time consuming activity. It takes hours to yield dynamic subsumption relations. Moreover, we have to rely on good automatic test suite generators. However, tests only improve confidence in the previous results since we do not have a proof [29].

In the approach presented here, we only need to encode the program and mutants (see Listing 8) to prove subsumption relations in few seconds. Analyzing the values given by the Z3 theorem prover for invalid theorems can help in this process to better understand why a mutation does not subsume another one.

By using our approach, we find some differences in the dynamic mutant subsumption graphs derived by Guimarães et al. [11] that contain integer expressions. All mutation subsumption relation graphs are different. We find that we cannot reduce the number of mutations for targets containing integer numbers. Since Guimarães et al. [11] rely on the test suite generators that do not consider all integer values, they find some subsumption relations different from our work.

All mutant subsumption relation graphs, proof scripts, and reproducibility instructions can be found in our notebook [8].

## 5. Evaluation

In our previous section, we analyze code fragments using weak mutation testing to derive a minimal set of mutants (see Table 1). This section evaluates our subsumption relations in the context of strong mutation testing by considering a complete program. Offutt and Lee [39] evaluated the effectiveness versus the efficiency of weak mutation testing. They found that weak mutation testing can be applied in a manner that is almost as effective as strong mutation testing and with significant computational savings. However the results using weak mutation testing do not always hold for strong mutation testing. For instance, Lindström and Márki [29] found that their subsumption relations for ROR identified using weak mutation testing do not hold for strong mutation testing.

To analyze to what extent our results hold for complete programs, first we change MUJAVA to include the results presented in Table 1 for 24 mutation targets. This tool is called MUJAVA-M [11]. Then we compare the results for MUJAVA and MUJAVA-M for a number of mutants generated from real projects in this section.

This section is organized as follows. First we present our research questions in Section 5.1. Section 5.2 presents the experimental planning. Section 5.3 explains the experimental procedure. Section 5.4 shows our results. We compare our technique to random sampling in Section 5.5. Finally, we discuss some threats to validity in Section 5.6. All data, setups, scripts, and MUJAVA-M are available in our companion website [8].

### 5.1. Research Questions

To better structure our evaluation, we rely on the Goal, Question, Metrics methodology [3]. The goal of our experiment consists of analyzing our approach,

implemented by MUJAVA-M, with the purpose of evaluating the subsumption relations we found in Z3 with respect to the number of mutants discarded (effort reduction), and the correctness of this reduction (effectiveness) from the point of view of testers in the context of applying mutation testing to Java open source programs (strong mutation testing).

To achieve this goal, we address the following research questions:

RQ<sub>1</sub>: *How many mutants are subsumed (effort reduction)?*

To answer this question, we count the number of mutants generated by MUJAVA and MUJAVA-M for each mutation target. Notice that answering RQ<sub>1</sub> is important because it allows us to estimate the amount of computational effort saved. The subsumption relations we embedded in MUJAVA-M must be effective. They should not discard important mutants that would be in a minimal set. To better understand this point, we formulate the following complementary research question:

RQ<sub>2</sub>: *How many mutants are incorrectly discarded from a minimal set (effectiveness)?*

To answer RQ<sub>2</sub>, we rely on the definition of minimal test set [1]. According to Amman et al. [1], a minimal test set necessary to kill a minimal mutants set must also kill all the mutants in the full mutants set. Thus, we generate this minimal test set and execute against the full mutants set. If a mutant from the full mutants set survives, this means that we incorrectly discarded this mutant. We compute the frequency of these cases.

## 5.2. Planning

We use five large open source programs to carry out our evaluation. Table 3 illustrates the studied programs, i.e., *joda-time*, *commons-math*, *commons-lang*, *h2*, and *javassist*. These programs vary in size and application domain. *joda-time* is a time manipulation library. *commons-math* is a library of mathematics and statistics components. *commons-lang* is a package of Java utility classes for

the classes that are in `java.lang`’s hierarchy. *h2* is a Java SQL-based database. *javassist* is class library for editing bytecodes. We performed the evaluation on Intel Core i5-7400 with 8 GB of RAM equipped with Linux 3.10.0 operating system. We used MUJAVA and MUJAVA-M command-line version. In both cases, all method-level mutation operators were enabled.

Table 3: Programs used in our evaluation.

Project	Version	Lines of Code (LOC)
joda-time	2.10.1	28,790
commons-math	3.6.1	100,364
commons-lang	3.6	27,267
h2	1.4.199	134,234
javassist	3.20	35,249

After generating mutants with each tool, we need to calculate the incorrectly discarded mutants by MUJAVA-M. Thus, we need to execute a minimal test set — necessary to kill the MUJAVA-M mutants — against the mutants generated by MUJAVA. To find out the minimal test set, we rely on EVOSUITE’s [7] Regression test suite generation (EVOSUITER) version 1.0.6. EVOSUITER is a specialization of EVOSUITE that tries to generate one test revealing the difference between two versions of a Java class. For instance, given two Java classes with a small syntactic difference in code, say a mutant, EVOSUITER tries to find a test case that exposes this behavioral difference between the two files. We set up 60 seconds as the time limit to EVOSUITER generate tests. We used the default values for the other parameters.

In case the mutant survives the test generated by EVOSUITER, we try to discard equivalent mutant. Equivalent mutants contribute negatively to the confidence assessment of the reduction applied. Unfortunately, detecting equivalent mutants is a well-known undecidable problem [4]. To minimize this problem, we avoid some equivalent mutants by using the Trivial Compiler Equivalence (*TCE*) [25]. *TCE* is a sound tool because it checks whether the bytecodes of the

original program and the mutant are the same. This eliminates the possibility of false positives. However, *TCE* cannot identify equivalent mutants that have different bytecodes, which may yield false negatives.

In summary, to answer RQ<sub>1</sub> and RQ<sub>2</sub> the plan is the following: generate the mutants with MUJAVA and MUJAVA-M, then generate the minimal tests set with EVOSUITER, execute the test generated against the MUJAVA mutants, detect equivalence with *TCE*, and calculate the surviving mutants. Because it is a computationally costly experiment, we leave the programs running for seven days for each subject. Consequently, the number of randomly selected files of each subject varied. In total, we evaluate 125 class files (see Table 4).

### 5.3. Procedure

We explain how we proceed to answer the research questions RQ<sub>1</sub> and RQ<sub>2</sub>. A Java class is the MUJAVA unity of work, thus we need to generate the mutants for the whole class. Applying all possible mutants to all files in a large program is clearly infeasible. This way we randomly selected a set of Java class files for each subject. With the classes selected, we executed MUJAVA and MUJAVA-M against these classes to generate the full set and a minimal set of mutants, respectively. We enabled all method-level mutation operators in both tools.

Next, we added the mutants of MUJAVA and MUJAVA-M grouped by target. For instance, for each target  $t$  in a given class file, MUJAVA generated the full set  $M = \{m_1, m_2, m_3, m_4\}$  containing all mutants, and MUJAVA-M generated a minimal set  $\bar{M} = \{m_1, m_2\}$  containing only the sufficient mutants according to the subsumption relations found previously by our approach (Section 4).

We now proceed to create a minimal test set. As explained, a minimal test set necessary to kill the minimal mutants set must also kill the full mutants set [1]. Thus, we use EVOSUITER to create a test case for each mutant in a minimal set ( $\bar{M}$ ). We provide the original program and a mutant from  $\bar{M}$ , and EVOSUITER generates a test containing only one test case to kill the mutant. We repeat this process for all mutants in  $\bar{M}$ . At the end, we group the generated tests to create a minimal test suite  $\bar{T}$  for a minimal set of mutants  $\bar{M}$ .



To validate if the mutants of  $\bar{M}$  indeed represent a minimal mutants set for the target, we execute  $\bar{T}$  against  $M$ . In case all mutants of  $M$  get killed, we confirm that  $\bar{M}$  is a reliable representation of  $M$ . But if a mutant of  $M$  survives, it represents a fail in our approach. For example, if only the  $m_1$ ,  $m_2$ , and  $m_3$  mutants of  $M$  are killed by suite  $\bar{T}$ , only 75% of the mutants in the full set were killed. This means that  $m_4$  is a useful mutant and should not be discarded from a minimal set. An exception occurs when  $m_4$  is an equivalent mutant. In this case  $m_4$  is useless to the mutation test. This way, we executed *TCE* against the mutants of  $M$  that survived to  $\bar{T}$ . If *TCE* identifies a mutant as equivalent, we take this mutant out of the analysis. If *TCE* does not mark a mutant as equivalent, then we understand that this mutant represents an error in our reduction and it should be part of the minimal mutant set.

To understand if our approach has eliminated important mutants, we verified the number of mutants not generated by MUJAVA-M that should be part of a minimal set. We also manually verified a subset of these incorrectly deleted mutants.

To automate the process described before, we create a script that executes all steps. In some exceptional scenarios we discard the target. Below we list these scenarios:

- If EVOSUITER cannot identify a test case to distinguish the original program and a mutant in a limit of 60 seconds, we did not proceed with the analysis of the target.
- We execute the minimal test suite against the original program to confirm they are passing. We repeat this process three times to reduce the presence of flaky tests [30]. In case we identify flaky tests, or the test suite does not pass in the original program, we do not proceed with the analysis of the target.

#### 5.4. Results

Next we answer our research questions.

#### 5.4.1. RQ<sub>1</sub>: How many mutants are subsumed (effort reduction)?

Table 4 presents the number of mutants generated by MUJAVA and MUJAVA-M for each subject. In particular, we analyzed 1,403 occurrences of mutation targets in 125 classes. MUJAVA generated 6,898 mutants, which gives an average of 4.92 mutants per target. MUJAVA-M, in its turn, generated 1,850 mutants for the same set of mutation targets, i.e., an average of 1.32 mutants per target. This way, MUJAVA-M achieved an average reduction of 71.38% in the number of generated mutants when compared to the original version of MUJAVA.

Table 4: Number of mutants per subject.

Project	Classes	MuJava	MuJava-M
joda-time	38	2,755	666
commons-math	34	1,282	368
commons-lang	22	1,737	537
h2	11	231	63
javassist	20	893	216
<b>Total</b>	125	6,898	1,850

Table 5 illustrates the occurrences of 18 mutation targets we analyzed in the 125 classes. The most common target is `exp`. We identified 1,089 `exp` occurrences. The effort reduction rate is 75.80% on average for this target, respectively.

Table 5: General results for some targets.

Mutation Target	Occurrences	Projects	Reduction	Effectiveness
lexp > rexp	19	4	62.50%	100.00%
lexp >= rexp	26	4	62.50%	100.00%
lexp < rexp	35	4	62.50%	100.00%
lexp <= rexp	16	3	62.20%	100.00%
lexp == rexp	34	2	62.50%	100.00%
lexp != rexp	14	4	62.50%	100.00%
lexp && rexp	13	2	55.56%	100.00%
lexp    rexp	25	4	55.56%	100.00%
lexp & rexp	33	2	64.32%	100.00%
lexp   rexp	6	1	40.00%	100.00%
lexp ^ rexp	6	1	83.33%	100.00%
exp	1,089	5	75.80%	47.20%
!exp	27	4	50.00%	100.00%
-exp	38	4	58.70%	92.68%
~exp	13	2	58.06%	100.00%
exp++	4	3	71.43%	100.00%
exp--	4	2	83.33%	66.67%
lhs ^= rhs	1	1	66.67%	50.00%

We achieve significant reductions when considering the total number of generated mutants (see column “Reduction” in Table 5). However, we may have discarded important mutants for the mutation analysis. In this sense, to better understand to what extent our reductions are indeed focusing only on redundant mutants, we now answer RQ<sub>2</sub>.

*5.4.2. RQ<sub>2</sub>: How many mutants are incorrectly discarded from a minimal set (effectiveness)?*

Table 5 also presents numbers with respect to the effectiveness of MUJAVAM, i.e., we check whether the mutants discarded by our tool were indeed dis-

carded correctly. Column “Effectiveness” presents these results. This percentage represents the number of mutants generated by MUJAVA that were killed by the minimal test set. In the ideal scenario, the minimal test set should kill all MUJAVA mutants.

According to Table 5, we achieve 100% of effectiveness in 14 targets. On the other hand, we achieved only 47.20% of effectiveness for the `exp` target (the one more common in the subjects we studied, i.e., 1,089 occurrences). Notice that `exp` is a very generic constructor that can be, for example, a variable that stores the index of an array, i.e., `arr[i]`.

There are some reasons why the results presented in Table 1 for mutation targets in isolation do not hold for some projects (Table 5). There are scenarios in which we infect the program state but the infection is not propagated [56]. So we cannot kill the mutants. The changes are not observable for users. It is an internal change. This is a limitation of weak mutation testing [13]. However, we also have limitations in the procedure presented in Section 5.3. There are some challenges in using automatic test suite generators [54, 55]. These challenges negatively impact on the effectiveness of our approach. For instance, they may not generate some input values to kill some mutants. So, the test suite generator cannot infect the program state [56]. There are some limitations in the automatic test suite generators related to defining oracles [54]. We have scenarios in which the program state is infected, the infection is propagated, but we cannot reveal it since we do not have a good oracle. The automatic test suite generators may generate flaky (unstable) tests [54]. As future work, we intend to manually analyze our sample, and also consider projects with test suites created by developers.

Next, we discuss an error when defining a minimal set for an instance of the target `lhs ^= rhs` [11]. This target occurred only once in the subjects studied (see Table 5). Listing 12 presents a code snippet of the `BooleanUtils` class of the project *commons-lang*. At Line 5 of the `xor` method there is the following statement: `result ^= element`. This statement applies an exclusive disjunction logic operation among all elements of the array. The minimal mutation

set for this target is made up of just one mutation: `ASRS(=)` as presented in Table 1. However, the minimal test set did not kill the `ASRS(&=)` mutation.

Listing 12: Code snippet from *commons-lang* project.

```
public static boolean xor(final boolean... array) {
    ...
    boolean result = false;
    for (final boolean element : array) {
        result ^= element;
    }
    return result;
}
```

Lindström and Márki [29] suggested that the subsumption relations cannot hold when the mutated statements are re-executed (in the context of strong mutation [6]). If the mutated instruction is executed more than once by any test execution, we cannot determine the future state of the program. Notice that the mutation target is inside the for loop (see Listing 12). Since our subsumption relations were obtained using weak mutation testing, they are not sufficient to represent the mutation within a repeating context.

In summary, we show that our approach to identify subsumption relations in Z3 using weak mutation testing (see Section 4) has a good balance between the effort (sampling rate of 28.62%) required to derive them and the effectiveness (75.93%) for the targets considered in our evaluation in the context of strong mutation testing.

### 5.5. Random Sampling

Gopinath et al. [9] compared the effectiveness of some mutation reduction strategies to random sampling. In their evaluation, none of the mutation reduction strategies evaluated produced an effectiveness advantage larger than 5% in comparison with random sampling. In summary, they argue that mutation reduction strategies are considered harmful. In this section, we compare our approach to random sampling.

We consider the programs presented in Table 3 and targets presented in Table 5 to evaluate the random sampling strategy. The baseline minimal set of mutants is defined by joining a minimal set of mutants identified by our approach, and the set of nonequivalent mutants not killed by our approach. For each target, we randomly select the mutant set and count the number of mutants in the baseline minimal set of mutants. To avoid bias, we repeat this process 100 times and yield the median value.

We use 10 sampling rates from 0 to 100%. Figure 5 presents our results. The random sampling approach yields an average effectiveness of 75% in correctly identifying the baseline minimal set of mutants when it uses a sampling rate of 60%. Our approach presented in Section 4 yields an effectiveness of 75.93% when it uses a sampling rate of 28.62%. This way, different from the results obtained by Gopinath et al. [9] in their study, using the random sampling strategy is not a good approach compared to ours. Our reduction mutation strategy is not considered harmful.

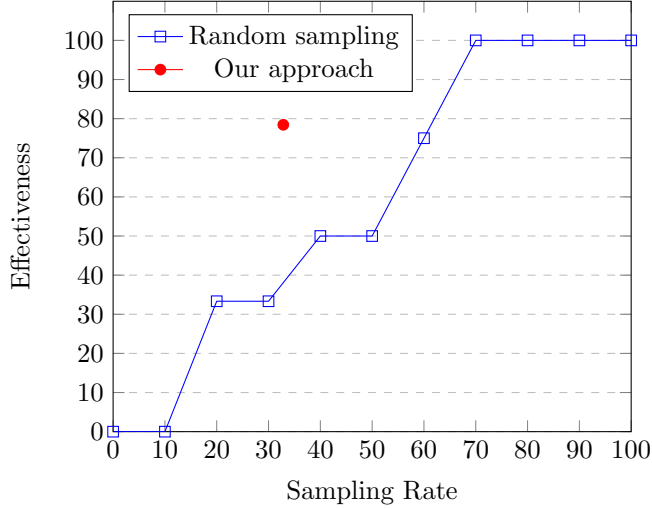


Figure 5: Comparing our approach to random sampling strategy.

### 5.6. Threats to Validity

The set of projects we used represents a threat to external validity. Also, we did not evaluate all files of all projects. To increase diversity, we consider projects of different sizes and domains. As another threat to external validity, we focused only on method-level operators of only one tool, i.e., MUJAVA for Java. In some cases MUJAVA generates mutants that do not compile or fails to generate some mutants, representing a threat to internal validity.

We only considered in this study mutation targets that did not generate flaky tests and that EVOSUITER could generate the minimal test sets. This represents a threat to internal validity. This decision was necessary to assess the effectiveness of the reductions. The minimal test sets also poses a threat to internal validity. This is because computing minimal mutant sets for all possible test sets is computationally hard [1]. Thus, the EVOSUITER can generate the test set which is minimal but not minimum [1].

The mutants that survived the minimal test sets also represents a threat. Despite running *TCE* to identify equivalent mutants, *TCE* cannot detect all equivalent mutants due to the undecidability of the Equivalent Mutant Problem [34].

Some targets did not appear frequently in our evaluation. For instance, the mutation target `lhs ^= rhs`, occurred only once. So, the effectiveness of a minimal set defined for some targets may not hold for general cases. We intend to perform other studies to evaluate these targets.

## 6. Related Work

There are some strategies to reduce costs for mutation analysis in the literature [50]. Kaminski et al. [24] defined the mutant subsumption graphs for six targets: `lexp > rexp`, `lexp >= rexp`, `lexp < rexp`, `lexp <= rexp`, `lexp == rexp`, and `lexp != rexp`. We yield the same minimal set for them. Moreover, we encode more targets, presented in Table 1, using the Z3 theorem prover. Using a similar strategy, Just et al. [21] presented sufficient sets of non-redundant

mutations for the `COR` and `UOI` operators. These subsumption hierarchies are defined by manually analyzing the combinations of all possible input situations. However, in several other cases, analyzing all possible combinations is prohibitive due to the high costs. Our approach encodes a theory in `Z3` and uses the `Z3` theorem prover to automatically deduce the subsumption relations.

Guimarães et al. [11] proposed an approach to identify subsumption relations using automatic test suite generators in the context of strong mutation testing. In contrast, we propose an approach that is simpler to derive subsumption relations. Indeed, we do not need to generate and compile a number of mutants. We do not need to automatically generate tests, nor execute them. Instead by using our theory, we have to encode the program and mutation operators. Then the `Z3` theorem prover automatically proved a number of subsumption relations for weak mutation testing.

Just and Schweiggert [23] presented a study that analyzes the effect of redundant mutants on mutation analysis efficiency, mutation score, and mutation coverage ratio. They show that the mutants generated by `COR`, `ROR`, and `UOI` have a mean ratio of 45% of the total mutants generated. Using the sufficient set of non-redundant mutations for these operators, the number of mutants was reduced by 27% overall. Just and Schweiggert also show that redundant mutants worsen the accuracy of the mutation score.

Papadakis and Malevris [47] showed that random selection of subsets containing 10%-60% of the generated mutants reduces the ability to detect failures by 26%-6%, respectively. Offutt et al. [38] presented an empirical approach to define an appropriate set of selective mutation operators. The idea was to randomly select a subset of mutation operators [35, 57]. Perez et al. [5] explored Evolutionary Mutation Testing to reduce the number of mutants to be executed. Namin et al. [37] formulated the selective mutation problem as a statistical problem. They applied linear statistical approaches to identify a subset of 28 mutation operators for `C`. Some techniques used clustering algorithms to reduce the number of mutants by selecting only a subset of mutants from each cluster [16, 14]. Other strategies [17, 28] for reducing costs uses the idea of



higher order mutants (mutants with more than one syntactic change), which subsume the behavior of two or more mutants with only one syntactic change, also known as first order mutants. We show how to encode two higher order mutants in our approach.

However, in another study, Gopinath et al. [9] found no differences in effectiveness between selective mutation and random selection. The main challenge in reducing the mutants set is not losing useful information. We show that our approach has a better effectiveness than the random sampling strategy for the same sampling rate. Just et al. [22] stated that existing approaches to selective mutation do not take program context into account, and this is fundamental to avoid losing useful information.

The high cost of mutation testing creates an entry barrier to its use in the software industry, but the effectiveness of mutation testing in assessing the quality of the test suite makes it attractive. Therefore, there is an incentive to carry out cost-saving studies and alternative ways to use mutation, such as the approach used by Google, where only one mutant per target is chosen by a software engineer manually during the code quality inspection [48].

In our work, we propose to use subsumption relationships to reduce costs for mutation testing. Our approach is related to the selective mutation strategy, as we use the subsumption relationships found to select the most representative mutants among all generated mutants. Moreover, we encode a theory of subsumption relations in Z3, and use its theorem prover to identify a number of subsumption relations. We focus on identifying subsumption relations using weak mutation testing. We have to be careful when leveraging our results for strong mutation testing. Lindström and Márki [29] studied the subsumption relations between ROR mutants. They showed that ROR fault hierarchies identified using weak mutation testing do not hold for strong mutation testing. The problem may be mitigated by avoiding loop structures. We evaluate our approach in the context of strong mutation testing in Section 5, and show that our approach has a good balance between the effort required to derive the mutant subsumption relations and the effectiveness for the targets considered in

our evaluation.

Previous approaches focused on proposing approaches to detect equivalent mutants [19, 34]. Baldwin and Sayward used compiler optimizations [2] to detect equivalent mutants by checking whether the original program and the optimized program are identical. Kintis et al. [25] proposed the Trivial Compiler Equivalence for C and Java, and mutation tools (MILU and MUJAVA).

Offut and Pan [41, 42] developed a technique to detect equivalent mutants based on mathematical constraints that introduce a set of strategies to formulate the killing conditions of the mutants. If these conditions are not feasible, the mutant is equivalent. Voas and McGraw [56] and Hierons et al. [12] suggested to use program slicing to help with equivalence identification. These approaches suffer from inherent limitations in the scalability of constraint handling and slicing technology. Grun et al. [10] and Shuller and Zeller [52, 53] proposed that changes in coverage can be used to detect non-equivalents mutants. Shuler et al. [51] used invariants violation as a way to classify killable mutants. In our approach, we defined the function `keepNonEquivalentMutants` and used the Z3 theorem prover to identify equivalent mutants using weak mutation testing.

## 7. Conclusion

In this work, we automatically identify and prove a number of subsumption relations for method-level mutations using the Z3 theorem prover. Developers only need to specify the types and mutations in our encoding to identify subsumption relations (see Listing 8). In few seconds, the Z3 theorem prover automatically proves a number of subsumption relations for 37 mutation targets. We reduce the number of mutations in a number of mutation targets containing integer and boolean expressions. We show some examples on how to encode them and identify subsumption relations. We can extend our theory to consider other types of expressions. To evaluate our approach, we extend MUJAVA with some of our results and evaluate it in 125 classes of 5 real projects. Our tool achieves an effectiveness of 75.93% and a sampling rate of 28.62% in the context

of strong mutation testing. Moreover, we show that our approach is better than the random sampling strategy.

The results may help to build better mutation testing tools that will allow to reduce the mutation testing costs. We recommend the community to follow a similar approach presented here before proposing new mutations. We must propose new mutations that subsume the previous ones. In this way, developers can use a minimal set of mutations, hence reducing mutation testing costs. Overall, our work leverages lightweight formal methods to mutant analysis, resulting in effective gains for developers.

As future work, we intend to prove more subsumption relations by considering real numbers, other language constructs and mutations, and encoding more higher order mutants. It will require to encode the Java semantics of some constructions, such as classes and fields. We may encode the Java Featherweight semantics in Z3 [15] or in other systems in which we can interactively perform the proofs, such as PVS [43]. In this case, we may address some proofs that cannot be done (when `prove` yields `unknown`) using the Z3 theorem prover. Finally, the expressions and commands considered in this work for Java have a similar semantics in other languages, such as Python and C#. We intend to check whether the subsumption relations found also hold for other languages in the context.

## Acknowledgments

We would like to thank the anonymous reviewers for their insightful suggestions. This work was partially supported by CNPq and CAPES grants.

## References

- [1] Ammann, P., Delamaro, M.E., Offutt, J., et al., 2014. Establishing theoretical minimal sets of mutants, in: Proceedings of the International Conference on Software Testing, Verification, and Validation, IEEE. pp. 21–30.

- [2] Baldwin, D., Sayward, F., 1979. Heuristics for Determining Equivalence of Program Mutations. Technical Report. DTIC Document.
- [3] Basili, V.R., Caldiera, G., Rombach, H.D., 1994. The Goal Question Metric approach, in: Encyclopedia of Software Engineering. Wiley.
- [4] Budd, T., Angluin, D., 1982. Two notions of correctness and their relation to testing. *Acta Informatica* 18, 31–45.
- [5] Delgado-Pérez, P., Segura, S., Medina-Bulo, I., 2017. Assessment of C++ object-oriented mutation operators: A selective mutation approach. *Software Testing, Verification and Reliability* 27.
- [6] DeMillo, R.A., Lipton, R.J., Sayward, F.G., 1978. Hints on test data selection: Help for the practicing programmer. *Computer* 11, 34–41.
- [7] Fraser, G., Arcuri, A., 2011. EvoSuite: automatic test suite generation for object-oriented software, in: Proceedings of the Foundations of Software Engineering, ACM. pp. 416–419.
- [8] Gheyi, R., Ribeiro, M., Sousa, B., Guimarães, M., d’Amorim, M., Alves, V., Teixeira, L., Fonseca, B., 2020. Identifying method-level mutant subsumption relations using Z3 (artifacts). At [https://colab.research.google.com/drive/1VPzo0NEt8F\\_1XN8pioKy5Vd1Ip2fH9WS](https://colab.research.google.com/drive/1VPzo0NEt8F_1XN8pioKy5Vd1Ip2fH9WS).
- [9] Gopinath, R., Ahmed, I., Alipour, M., Jensen, C., Groce, A., 2017. Mutation reduction strategies considered harmful. *IEEE Transactions on Reliability* 66, 854–874.
- [10] Grün, B.J., Schuler, D., Zeller, A., 2009. The impact of equivalent mutants, in: Proceedings of the International Conference on Software Testing, Verification and Validation Workshops, pp. 192–199.
- [11] Guimarães, M., Fernandes, L., Ribeiro, M., d’Amorim, M., Gheyi, R., 2020. Optimizing mutation testing by discovering dynamic mutant subsumption

- relations, in: Proceedings of the International Conference on Software Testing, IEEE. pp. 98–208.
- [12] Hierons, R., Harman, M., Danicic, S., 1999. Using program slicing to assist in the detection of equivalent mutants. *Software Testing, Verification and Reliability* 9, 233–262.
  - [13] Howden, W., 1982. Weak mutation testing and completeness of test sets. *IEEE Transactions on Software Engineering* 8, 371–379.
  - [14] Hussain, S., 2008. Mutation clustering. Master’s thesis. Kings College London.
  - [15] Igarashi, A., Pierce, B.C., Wadler, P., 2001. Featherweight Java: A minimal core calculus for Java and GJ. *ACM Transactions on Programming Languages and Systems* 23, 396–450.
  - [16] Ji, C., Chen, Z., Xu, B., Zhao, Z., 2009. A novel method of mutation clustering based on domain analysis, in: Proceedings of the International Conference on Software Engineering and Knowledge Engineering, pp. 422–425.
  - [17] Jia, Y., Harman, M., 2008. Constructing subtle faults using higher order mutation testing, in: Proceedings of the International Working Conference on Source Code Analysis and Manipulation, IEEE. pp. 249–258.
  - [18] Jia, Y., Harman, M., 2009. Higher order mutation testing. *Information and Software Technology* 51, 1379–1393.
  - [19] Jia, Y., Harman, M., 2011. An analysis and survey of the development of mutation testing. *IEEE Transactions on Software Engineering* 37, 649–678.
  - [20] Just, R., Jalali, D., Inozemtseva, L., Ernst, M.D., Holmes, R., Fraser, G., 2014. Are mutants a valid substitute for real faults in software testing?, in: Proceedings of the Foundations of Software Engineering, pp. 654–665.

- [21] Just, R., Kapfhammer, G.M., Schweiggert, F., 2012. Do redundant mutants affect the effectiveness and efficiency of mutation analysis?, in: Proceedings of the International Conference on Software Testing, Verification and Validation, IEEE. pp. 720–725.
- [22] Just, R., Kurtz, B., Ammann, P., 2017. Inferring mutant utility from program context, in: Proceedings of the International Symposium on Software Testing and Analysis, pp. 284–294.
- [23] Just, R., Schweiggert, F., 2015. Higher accuracy and lower run time: efficient mutation analysis using non-redundant mutation operators. *Software Testing, Verification and Reliability* 25, 490–507.
- [24] Kaminski, G., Ammann, P., Offutt, J., 2011. Better predicate testing, in: Proceedings of the International Workshop on Automation of Software Test, ACM. pp. 57–63.
- [25] Kintis, M., Papadakis, M., Jia, Y., Malevris, N., Traon, Y.L., Harman, M., 2017. Detecting trivial mutant equivalences via compiler optimisations. *IEEE Transactions on Software Engineering* 44, 308–333.
- [26] Kintis, M., Papadakis, M., Malevris, N., 2010. Evaluating mutation testing alternatives: A collateral experiment, in: Proceedings of the Asia Pacific Software Engineering Conference, IEEE. pp. 300–309.
- [27] Kurtz, B., Ammann, P., Delamaro, M.E., Offutt, J., Deng, L., 2014. Mutant subsumption graphs, in: Proceedings of the International Conference on Software Testing, Verification and Validation workshops, IEEE. pp. 176–185.
- [28] Langdon, W.B., Harman, M., Jia, Y., 2009. Multi objective higher order mutation testing with genetic programming, in: Testing: Academic and Industrial Conference-Practice and Research Techniques, IEEE. pp. 21–29.

- [29] Lindström, B., Márki, A., 2019. On strong mutation and the theory of subsuming logic-based mutants. *Software Testing, Verification and Reliability* 29, e1667.
- [30] Luo, Q., Hariri, F., Eloussi, L., Marinov, D., 2014. An empirical analysis of flaky tests, in: *Proceedings of the Foundations of Software Engineering*, pp. 643–653.
- [31] Ma, Y.S., Offutt, J., . Description of MuJava’s method-level mutation operators. At <https://cs.gmu.edu/~offutt/mujava/mutopsMethod.pdf>.
- [32] Ma, Y.S., Offutt, J., Kwon, Y.R., 2005. MuJava: an automated class mutation system. *Software Testing, Verification and Reliability* 15, 97–133.
- [33] Ma, Y.S., Offutt, J., Kwon, Y.R., 2006. MuJava: a mutation system for Java, in: *Proceedings of the International Conference on Software Engineering*, ACM. pp. 827–830.
- [34] Madeyski, L., Orzeszyna, W., Torkar, R., Jozala, M., 2014. Overcoming the equivalent mutant problem: A systematic literature review and a comparative experiment of second order mutation. *IEEE Transactions on Software Engineering* 40, 23–42.
- [35] Mathur, A., 1991. Performance, effectiveness, and reliability issues in software testing, in: *Proceedings of the Annual International Computer Software and Applications Conference*, IEEE. pp. 604–605.
- [36] de Moura, L.M., Bjørner, N., 2008. Z3: an efficient SMT solver, in: *Proceedings of the Tools and Algorithms for the Construction and Analysis of Systems*, pp. 337–340.
- [37] Namin, A., Andrews, J., Murdoch, D., 2008. Sufficient mutation operators for measuring test effectiveness, in: *Proceedings of the International Conference on Software Engineering*, ACM. pp. 351–360.

- [38] Offutt, A.J., Lee, A., Rothermel, G., Untch, R.H., Zapf, C., 1996. An experimental determination of sufficient mutant operators. *ACM Transactions on Software Engineering and Methodology* 5, 99–118.
- [39] Offutt, A.J., Lee, S.D., 1994. An empirical evaluation of weak mutation. *IEEE Transactions on Software Engineering* 20, 337–344.
- [40] Offutt, J., 2011. A mutation carol: Past, present and future. *Information and Software Technology* 53, 1098 – 1107.
- [41] Offutt, J., Pan, J., 1996. Detecting equivalent mutants and the feasible path problem, in: *Proceedings of the Annual Conference on Computer Assurance*, pp. 224–236.
- [42] Offutt, J., Pan, J., 1997. Automatically detecting equivalent mutants and infeasible paths. *Software Testing, Verification and Reliability* 7, 165–192.
- [43] Owre, S., Rushby, J., Shankar, N., Stringer-Calvert, D., 1998. PVS: an experience report, in: *Applied Formal Methods—FM-Trends 98*, Springer-Verlag, Germany. pp. 338–345.
- [44] Pacheco, C., Lahiri, S., Ernst, M., Ball, T., 2007. Feedback-directed random test generation, in: *Proceedings of the International Conference on Software Engineering*, IEEE. pp. 75–84.
- [45] Papadakis, M., Henard, C., Harman, M., Jia, Y., Le Traon, Y., 2016. Threats to the validity of mutation-based test assessment, in: *Proceedings of the International Symposium on Software Testing and Analysis*, ACM. pp. 354–365.
- [46] Papadakis, M., Kintis, M., Zhang, J., Jia, Y., Traon, Y.L., Harman, M., 2019. Chapter six - mutation testing advances: An analysis and survey. *Advances in Computers* 112, 275–378.
- [47] Papadakis, M., Malevris, N., 2010. An empirical evaluation of the first and second order mutation testing strategies, in: *Proceedings of the 3rd*



International Conference on Software Testing, Verification, and Validation workshop, IEEE. pp. 90–99.

- [48] Petrovic, G., Ivankovic, M., 2018. State of mutation testing at Google, in: Proceedings of the International Conference on Software Engineering—Software Engineering in Practice, pp. 163–171.
- [49] Petrovic, G., Ivankovic, M., Kurtz, B., Ammann, P., Just, R., 2018. An industrial application of mutation testing: Lessons, challenges, and research directions, in: Proceedings of the International Conference on Software Testing, Verification and Validation Workshops, pp. 47–53.
- [50] Pizzoleto, A., Ferrari, F., Offutt, J., Fernandes, L., Ribeiro, M., 2019. A systematic literature review of techniques and metrics to reduce the cost of mutation testing. *Journal of Systems and Software* 157.
- [51] Schuler, D., Dallmeier, V., Zeller, A., 2009. Efficient mutation testing by checking invariant violations, in: Proceedings of the International Symposium on Software Testing and Analysis, pp. 69–80.
- [52] Schuler, D., Zeller, A., 2010. (un-)covering equivalent mutants, in: Proceedings of the International Conference on Software Testing, Verification and Validation, pp. 45–54.
- [53] Schuler, D., Zeller, A., 2013. Covering and uncovering equivalent mutants. *Software Testing, Verification and Reliability* 23, 353–374.
- [54] Shamshiri, S., Just, R., Rojas, J.M., Fraser, G., McMin, P., Arcuri, A., 2015. Do automatically generated unit tests find real faults? An empirical study of effectiveness and challenges, in: Proceedings of the Automated Software Engineering, pp. 201–211.
- [55] Soares, G., Gheyi, R., Murphy-Hill, E., Johnson, B., 2013. Comparing approaches to analyze refactoring activity on software repositories. *Journal of Systems and Software* 86, 1006 – 1022.

- [56] Voas, J., McGraw, G., 1997. Software fault injection: inoculating programs against errors. John Wiley & Sons, Inc.
- [57] Wong, W., Mathur, A., 1995. Reducing the cost of mutation testing: An empirical study. *Journal of Systems and Software* 31, 185–196.

### **CRedit Author Statement**

- Rohit Gheyi: Formal analysis, Methodology, Investigation, Conceptualization, Data Curation, Supervision, Writing - Original Draft
- Márcio Ribeiro: Methodology, Investigation, Conceptualization, Supervision, Writing - Original Draft
- Beatriz Sousa: Software, Data Curation, Validation, Writing - Review & Editing
- Marcio Guimarães: Software, Data Curation, Validation, Writing - Review & Editing
- Leo Fernandes: Software, Data Curation, Validation, Writing - Review & Editing
- Marcelo d'Amorim: Methodology, Investigation, Validation, Writing - Review & Editing
- Vander Alves: Formal analysis, Investigation, Writing - Review & Editing
- Leopoldo Teixeira: Formal analysis, Investigation, Writing - Review & Editing
- Balduino Fonseca: Methodology, Writing - Review & Editing

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: