

# Efficient Leveraging of Symbolic Execution to Advanced Coverage Criteria <sup>\*</sup>

Sébastien Bardin    Nikolai Kosmatov    François Cheynier  
CEA, LIST, Laboratoire pour la Sûreté du Logiciel  
91191, Gif-sur-Yvettes, France  
`first.name@cea.fr`

**Abstract**—Automatic test data generation (ATG) is a major topic in software engineering. In this paper, we bridge the gap between the coverage criteria supported by state-of-the-art white-box ATG technologies, especially Dynamic Symbolic Execution, and advanced coverage criteria found in the literature. We define a new testing criterion, label coverage, and prove it to be both expressive and amenable to efficient automation. We propose several innovative techniques resulting in an effective black-box support for label coverage, while a direct approach induces an exponential blow-up of the search space. Experiments show that our optimisations yield very significant savings, allowing to leverage ATG to label coverage with only a slight overhead.

**Keywords**—Testing, symbolic execution, coverage criteria

## I. INTRODUCTION

**Context and problem.** Automatic test data generation (ATG) is a major concern in software engineering and program analysis. Recent progress in automated theorem proving led to significant improvements of symbolic approaches for white-box ATG, such as Dynamic Symbolic Execution (DSE) [13], [36], [40]. Tools have been developed [5], [7], [8], [14], [37] and impressive case-studies have been carried out [7], [8], [16].

We consider the case where ATG aims at generating a test suite which is then passed to one or several external oracles, in order to assess for example functional correctness, security or performance. The more behaviours the test suite exercise, the better. A standard way of measuring this diversity passes through coverage criteria [2], [42]. Many such criteria have been defined along the years, from control-flow or data-flow criteria to mutation [11], input domain partitions and MCDC.

DSE mostly follows an exhaustive exploration of the path space of the program under test, covering all execution paths up to a given bound. While this path-oriented criterion proves successful in some contexts, it is well known that the resulting test suite can miss interesting behaviours related to data rather than control. Moreover, standard DSE does not support coverage objectives defined over artifacts not explicitly present in the source code, such as multiple-condition coverage, while they could efficiently guide test generation.

**Goal.** Our main objective is to bridge the gap between coverage criteria supported by symbolic ATG tools, especially DSE, and advanced coverage criteria found in the literature. Recent works aim at leveraging DSE to other coverage criteria [18],

[32], [33], [34], [35], [43] or improving DSE bug-detection abilities by making explicit run-time error conditions [9], [15], [20]. These approaches are mainly based on instrumentation and allow for black-box reuse of existing technologies. However, they come at a high price since they may induce a blow-up of the path space and a significant overhead (a recent paper [18, Table 2] reports on a 272x average time-overhead, with a worst case of 2,000x). We follow the same general line, emphasising black-box reuse as much as possible. However, we focus on two main points mostly left unaddressed: we want to formally characterize the class of coverage criteria that can be supported by DSE-like techniques, and we want to support it *efficiently*.

**Approach.** We define *label coverage*, a new testing criterion which appears to be both expressive and amenable to efficient automation. Especially, it turns out that *DSE can be extended for label coverage with only a slight overhead*. Labels are predicates attached to program instructions through a labelling function. A label is covered if a test execution reaches it and satisfies the predicate. This idea underlies former work [9], [15], [18], [20], [34], [43]. We generalize these results and propose ways of taming the potential blow-up. Especially, we introduce *tight instrumentation*, where “tight” is made precise in the paper, and a strong coupling of DSE and label coverage named *iterative label deletion*. Their combination results in an effective support for label coverage in DSE. Besides, both techniques can be implemented in black-box.

**Contribution.** Our main contributions are the following:

- We show that label coverage is expressive enough to faithfully emulate many standard coverage criteria, from decision or condition coverage (Theorem 1) to a substantial subset of weak mutations (the side-effect free fragment, Theorem 2). Labels can be seen in some way as a convenient and powerful specification mechanism for coverage criteria.
- We formally characterise the properties of direct instrumentation for label coverage. The instrumentation is sound w.r.t. label coverage and leads to very efficient coverage score computation. However, it yields an exponential increase as well as a “complexification” of the path space (Theorem 5).
- We propose *tight instrumentation* and *iterative label deletion* as ways of taming this complexity blow-up. Tight instrumentation yields only a linear growth of the path space without any complexification (Theorem 7). Both techniques are orthogonal and allow for

<sup>\*</sup> Work partially funded by EU FP7 (project STANCE, grant 317753) and French ANR (project BINSEC, grant ANR-12-INSE-0002).

a significant speed-up. Moreover, they can both be implemented either through dedicated DSE algorithms or in a black-box manner.

- We have implemented these results inside a DSE tool [40]. Initial experiments show that our optimisations yield very significant reductions of both search space and computation time compared to direct instrumentation (several-orders-of-magnitude speedup in some cases). It follows that ATG for label coverage can be achieved at a very reasonable cost w.r.t. usual DSE.

As a whole, label coverage forms the basis of a very generic and convenient framework for test automation, providing a powerful specification mechanism for test objectives and featuring efficient integration into symbolic ATG techniques as well as cheap coverage score computation. Moreover, static analysis techniques can also be used directly on the instrumented programs in order to detect uncoverable labels, as was proposed for mutation testing [24].

This work bridges part of the gap between symbolic ATG techniques and coverage criteria. On the one hand, we show that DSE techniques can be cheaply extended to more advanced testing criteria, such as side-effect free weak mutations. On the other hand, we identify a large subclass of weak mutations amenable to efficient automation, both in terms of ATG and mutation score computation.

**Outline.** After presenting basic notation (Section II), we define labels and explore their expressiveness (Section III). We then focus on automation. The direct instrumentation is defined and studied (Section IV). Afterwards, we describe our own approach to label-based ATG (Section V) and present first experiments (Section VI). Finally, we sketch a highly automated testing framework based on labels (Section VII), discuss related work (Section VIII) and give a conclusion (Section IX).

## II. BACKGROUND

### A. Notation

Given a program  $P$  over a vector  $V$  of  $m$  input variables taking values in a domain  $D \triangleq D_1 \times \dots \times D_m$ , a test datum  $t$  for  $P$  is a valuation of  $V$ , i.e.  $t \in D$ . The execution of  $P$  over  $t$ , denoted  $P(t)$ , is a path (or run)  $\sigma \triangleq (loc_1, S_1) \dots (loc_n, S_n)$ , where the  $loc_i$  denote control-locations (or simply locations) of  $P$  and the  $S_i$  denote the successive internal states of  $P$  ( $\approx$  valuation of all global and local variables as well as memory-allocated structures) before the execution of each  $loc_i$ . A test datum  $t$  reaches a location  $loc$  with internal state  $S$ , denoted  $t \rightsquigarrow_P (loc, S)$ , if  $P(t)$  is of the form  $\sigma_1 \cdot (loc, S) \cdot \sigma_2$ . A test suite  $TS$  is a finite set of test data.

Given a test objective  $\mathbf{c}$ , we write  $t \rightsquigarrow_P \mathbf{c}$  if test datum  $t$  covers  $\mathbf{c}$ . We extend the notation for a test suite  $TS$  and a set of test objectives  $\mathbf{C}$ , writing  $TS \rightsquigarrow_P \mathbf{C}$  when for any  $\mathbf{c} \in \mathbf{C}$ , there exists  $t \in TS$  such that  $t \rightsquigarrow_P \mathbf{c}$ . These definitions are generic and leave the exact definition of “covering” to the considered testing criterion. For example, test objectives derived from the Decision Coverage criterion are of the form  $\mathbf{c} \triangleq (loc, \text{cond})$  or  $\mathbf{c} \triangleq (loc, !\text{cond})$ , where  $\text{cond}$  is the condition of the branching instruction at location  $loc$ , and  $t \rightsquigarrow_P \mathbf{c}$  if  $t$  reaches some  $(loc, S)$  where  $\text{cond}$  evaluates to *true* (resp. *false*) in  $S$ .

Coverage criteria used through the paper and their associated notions of covering are described in Section III-B.

### B. DSE in brief

We recall here a few basic facts about Symbolic Execution (SE) [21] and Dynamic Symbolic Execution (DSE) [13], [36], [40]. Let us consider a program under test  $P$  with input variables  $V$  over domain  $D$  and a path  $\sigma$  of  $P$ . The key insight of SE is that it is possible in many cases to compute a *path predicate*  $\phi_\sigma$  for  $\sigma$  such that for any input valuation  $t \in D$ , we have:  $t$  satisfies  $\phi_\sigma$  iff  $P(t)$  covers  $\sigma$ . In practice, path predicates are often under-approximated and only the left-to-right implication holds, which is already fine for testing: SE outputs a set of pairs  $(t_i, \sigma_i)$  such that each  $t_i$  is ensured to cover the corresponding  $\sigma_i$ . Hence, SE is *sound* from a testing point of view. DSE enhances SE by interleaving concrete and symbolic executions. The dynamically collected information can help the symbolic step, for example by suggesting relevant approximations.

A simplified view of SE is depicted in Algorithm 1. While high-level, it is sufficient to understand the rest of the paper. We assume that the set of paths of  $P$ , denoted  $Paths(P)$ , is finite. In practice, DSE tools enforce this assumption through a bound on path lengths. We assume the availability of a procedure for path predicate computation (with predicates in some theory  $T$ ), as well as the availability of a solver taking a formula  $\phi \in T$  and returning either *sat* with a solution  $t$  or *unsat*. All DSE tools rely on such procedures. The algorithm builds iteratively a test suite  $TS$  by exploring all paths.

---

#### Algorithm 1: Symbolic Execution algorithm

---

**Input:** a program  $P$  with finite set of paths  $Paths(P)$   
**Output:**  $TS$ , a set of pairs  $(t, \sigma)$  such that  $P(t) \rightsquigarrow_P \sigma$

```

1  $TS := \emptyset;$ 
2  $S_{paths} := Paths(P);$ 
3 while  $S_{paths} \neq \emptyset$  do
4   | choose  $\sigma \in S_{paths}; S_{paths} := S_{paths} \setminus \{\sigma\};$ 
5   | compute path predicate  $\phi_\sigma$  for  $\sigma;$ 
6   | switch  $solve(\phi_\sigma)$  do
7     |   case  $sat(t): TS := TS \cup \{(t, \sigma)\};$ 
8     |   case  $unsat:$  skip ;
9   | endsw
10 end
11 return  $TS;$ 

```

---

The major issue here is that SE and DSE must in some ways explore all  $Paths(P)$ . Advanced tools explore this set lazily, yet they still have to crawl it. Therefore, the size of  $Paths(P)$ , denoted  $|Paths(P)|$ , is one of the two major bottlenecks of SE and DSE, the other one being the average cost of solving path predicates.

Note that Bounded model checking (BMC) [10] is sensitive to the same parameters, as it amounts to building a large formula encompassing all paths up to a given length.

## III. LABEL COVERAGE

### A. Definitions

Given a program  $P$ , a *label*  $l$  is a pair  $(loc, \varphi)$  where  $loc$  is a location of  $P$  and  $\varphi$  is a predicate such that:

- $\varphi$  contains only variables and expressions defined in  $P$  at location  $loc$ ;
- $\varphi$  contains no side-effect expressions.

An *annotated program* is a pair  $\langle P, L \rangle$  where  $L$  is a set of labels defined over  $P$ . A test datum  $t$  covers  $l \triangleq (loc, \varphi)$ , denoted  $t \rightsquigarrow_{\langle P, L \rangle} l$ , if  $t$  covers some  $(loc, S)$  with  $S$  satisfying predicate  $\varphi$ . We say that a test suite satisfies the *label coverage* testing criterion, denoted by **LC**, if it covers all labels in  $L$ .

For simplicity, we consider in the rest of the paper *normalized programs*, i.e. programs such that no side-effect occurs in any condition of a branching instruction. This is not a severe restriction since any (well-defined) program  $P_1$  can be rewritten into a normalized program  $P_2$ , using intermediate variables to evaluate the side-effect prone conditions outside the branching instruction. For example, `if (x++ <= y && e==f)` becomes `tmp = x++; if (tmp <= y && e==f)`. Notice that similar transformations are automatically performed by the Cil library [30] frequently used by DSE tools for C programs [36], [40].

### B. Expressiveness of label coverage

We seek to characterize the power of the **LC** testing criterion. A key notion is that of *labelling function*. A labelling function  $\psi$  maps a program  $P$  into an annotated program  $\psi(P) \triangleq \langle P, L \rangle$ .

*Definition 1:* A coverage criterion **C** can be simulated by **LC** if there exists a labelling function  $\psi$  such that for any program  $P$  and any test suite  $TS$ , we have  $TS \rightsquigarrow_P \mathbf{C}$  iff  $TS \rightsquigarrow_{\psi(P)} \mathbf{LC}$ .

We show first how **LC** can simulate basic graph and logic coverage criteria [2], [42]: instruction coverage **IC**, decision coverage **DC**, condition coverage **CC** (covering all atomic conditions appearing in each decision, and their negations), decision-condition coverage **DCC** (basically, **DC** plus **CC**) and multiple-condition coverage **MCC** (covering all combinations of atomic conditions). The basic idea is to introduce in  $P$  labels based on branching predicates and their atomic conditions. An example for **CC** is depicted in Figure 1, where additional labels (right) enforce coverage of the two atomic conditions  $x==y$  and  $a<b$ .

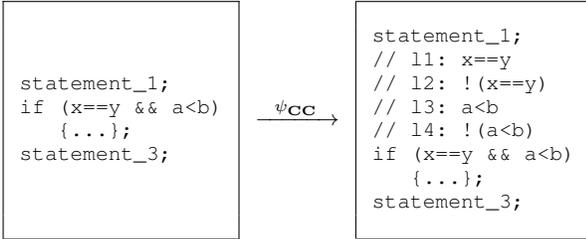


Fig. 1. Simulating **CC** with labels

*Theorem 1:* The coverage criteria **IC**, **DC**, **CC**, **DCC** and **MCC** can be simulated by **LC**.

*Proof:* We need to define a suitable labelling function for any of the considered criteria. For **IC**, we choose the labelling function  $\psi_{\mathbf{IC}}$  adding all labels of the form  $(loc, true)$ , where  $loc$  is any location of  $P$ . Given a test suite  $TS$ ,  $TS \rightsquigarrow_P \mathbf{IC}$  iff

$TS$  can reach any  $loc$  of  $P$  iff  $TS$  covers any  $(loc, true)$  iff  $TS \rightsquigarrow_{\psi_{\mathbf{IC}}(P)} \mathbf{LC}$ . We conclude that **IC** can be simulated by **LC**. Other criteria are handled similarly. The labelling function  $\psi_{\mathbf{DC}}$  adds the set of all  $(loc, \varphi)$  and  $(loc, \neg\varphi)$ , where  $loc$  contains a conditional statement with condition  $\varphi$ .  $\psi_{\mathbf{CC}}$  adds the set of all  $(loc, a_i)$  and  $(loc, \neg a_i)$ , where  $loc$  contains a conditional statement whose atomic conditions are exactly the  $a_i$ .  $\psi_{\mathbf{DCC}}$  adds the union of labels added by  $\psi_{\mathbf{DC}}$  and  $\psi_{\mathbf{CC}}$ .  $\psi_{\mathbf{MCC}}$  adds the set of all  $(loc, \bigwedge_i \bar{a}_i)$ , where the  $a_i$  are atomic conditions and  $\bar{a}_i$  denotes either  $a_i$  or  $\neg a_i$ . ■

**Weak mutations.** We now consider a more involved testing criterion, namely weak mutations. In mutation testing [11], test objectives consist of *mutants*, i.e. slight syntactic modifications of the program under test. In the strong mutation setting **M**, a mutant  $M$  is covered (or *killed*) by a test datum  $t$  if the output of  $P(t)$  differs from the output of  $M(t)$ . In the weak mutation setting **WM** [17], a mutant  $M$  is covered by  $t$ , denoted  $t \rightsquigarrow_P M$ , if the internal states of  $P(t)$  and  $M(t)$  differ from each other right after the mutated location (cf. Figure 2). **M** is a powerful testing criterion in practice [1], [28]. While less powerful in theory, **WM** appears to be almost equivalent to **M** in practice [26].

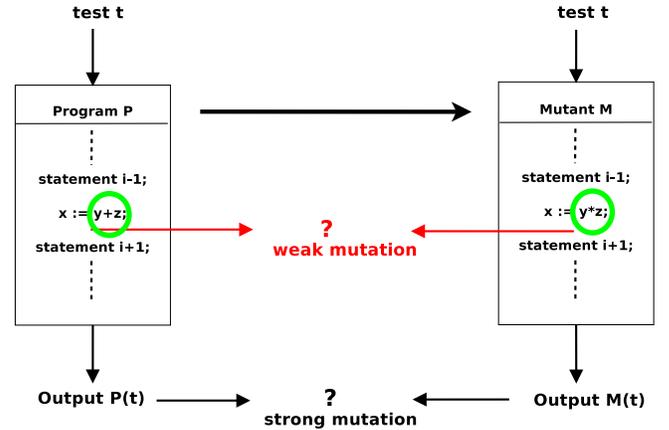


Fig. 2. Strong and weak mutations

We show hereafter that a substantial part of **WM** can be simulated by **LC**. First we need a few more definitions. Mutation testing is parametrised by a set of *mutation operators*  $O$ . A mutation operator  $op \in O$  is a function mapping a program  $P$  into a finite set of well-defined programs (mutants), such that  $P$  differs from each mutant  $M$  in only one location (atomic mutation). We denote  $\mathbf{WM}_O$  the weak mutation criterion restricted to mutants created through operators in  $O$ . We consider that mutations can affect either a left-hand side value (lhs), an expression or a condition. This is a very generic model of mutations, encompassing all standard operators [2]. Finally, we restrict ourselves to mutation operators neither affecting nor introducing side-effect expressions (including calls to side-effect prone functions). We refer to such operators as *side-effect free mutation operators*.

*Theorem 2:* For any finite set  $O$  of side-effect free mutation operators,  $\mathbf{WM}_O$  can be simulated by **LC**.

*Proof:* For simplicity, let us consider first a single mutation operator  $op \in O$ . The main idea is to introduce *one label*

for each mutant  $M$  created by  $op$ , so that covering the label is equivalent to distinguishing  $M$  from  $P$  once the modified location has been reached. This transformation is depicted in Figure 3. Let us consider a mutant  $M$  differing from  $P$  only at location  $loc$ . We consider three cases, depending on the modification introduced by  $op$ :

- $lhs := expr$  becomes  $lhs := expr'$ : we add label  $l \triangleq (loc, expr \neq expr')$ . We must prove that  $t \rightsquigarrow_P M$  iff  $t \rightsquigarrow_{\psi(P)} l$ . Note that  $t \rightsquigarrow_{\psi(P)} l$  iff  $t$  reaches  $loc$  with an internal state such that  $expr$  and  $expr'$  evaluate to different values. This is equivalent to say that  $P(t)$  and  $M(t)$  are in different internal states right after  $loc$ , which corresponds by definition to  $t \rightsquigarrow_P M$ .
- $if (cond)$  becomes  $if (cond')$ : we add label  $l \triangleq (loc, cond \oplus cond')$ , where  $\oplus$  is the xor-operator. We follow the same line of reasoning as in the previous case. The  $\oplus$  operator ensures that  $P(t)$  and  $M(t)$  will not follow the same branching condition.
- $lhs := expr$  becomes  $lhs' := expr$ : we add label  $l \triangleq (loc, \alpha(lhs) \neq \alpha(lhs') \wedge (lhs \neq expr \vee lhs' \neq expr))$ , where  $\alpha(x)$  denotes the memory location ( $\approx$  address) of  $x$ , not its value. For example, in C the memory location is given by the  $\&$  operator. This case requires a little bit more explanation. In order to observe a difference between  $P(t)$  and  $M(t)$  right after the mutated location, we need first that  $lhs'$  and  $lhs$  refer to different memory locations (which is not always obvious in the case of aliasing expressions). Moreover, there are only two ways of noticing a difference: either the old value of  $lhs$  differs from  $expr$ , then  $lhs$  will evaluate to different values in  $P(t)$  (equals to  $expr$ ) and in  $M(t)$  (remains unchanged) just after the mutation, or the symmetric counterpart for  $lhs'$ . This is exactly what  $l$  encodes.

By iterating this technique on all mutants created by the considered mutation operators, we obtain the desired labelling function. ■

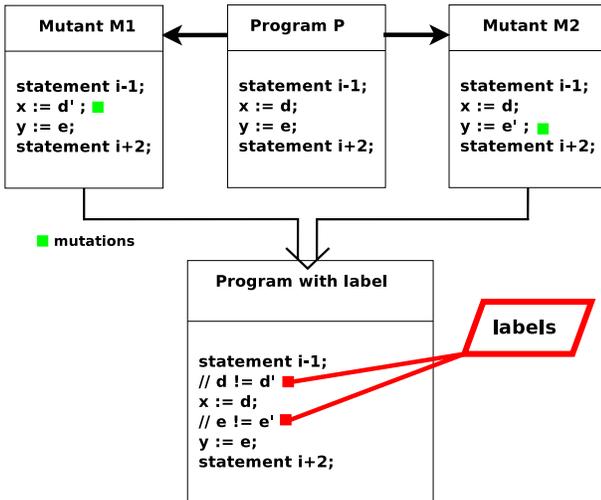


Fig. 3. Simulating weak mutants with labels

The subset of mutations we have been considering so far is limited to (1) atomic mutations and (2) side-effect free

operators. The first restriction is not a major issue as atomic mutations have been proved to be almost as powerful as high-order mutations [25]. The second restriction has two sides: (2.a) it forbids mutation operators *introducing* side-effects, for example mapping  $x$  to  $x++$ , and (2.b) it forbids to mutate a side-effect prone expression. Restriction (2.a) is not severe: it encompasses operators ABS, ROR, AOR, COR and UOI [2], which have been shown mostly equivalent to much larger sets of operators [27], [39]. It is left as an open question to quantify more precisely what is lost with restriction (2.b). Still, the previous points show that side-effect free **WM** is a substantial (strict) subset of **WM**.

**Other criteria.** Several other testing criteria commonly found in the literature can be emulated by labels. For example:

- **Input Domain Coverage:** assuming a partition of the input domain  $D$  of  $P$  given as disjoint predicates  $\varphi_1, \dots, \varphi_k$ , this criterion consists of considering one  $t_i$  for each  $\varphi_i$ . The corresponding labelling function adds all labels of the form  $(loc_0, \varphi_j)$ , where  $loc_0$  is the entry point of  $P$ . The approach is independent of the way the partition is obtained, covering both interface-based and functionality-based partitions [2].
- **Run-Time Error Coverage:** test objectives corresponding to *run-time errors* such as those implicitly searched for in active testing or assertion-based testing [9], [15], [20] can be easily captured by labels, including division by zero, out-of-bound array accesses or null-pointer dereference. Typically, any error-prone instruction at location  $loc$  with a precondition  $\varphi_{safe}$  will be tagged by a label  $(loc, \neg\varphi_{safe})$ .

**Limits.** The following criteria cannot be emulated through labels, at least with simple encoding: weak mutations with operators involving side-effects, criteria imposing constraints on paths rather than constraints on program locations (k-path coverage, data-flow criteria) and criteria involving different paths (MDCD) or even different programs (strong mutations). It is left as future work to study whether these limitations are strict or not.

#### IV. AUTOMATING LC: A FIRST ATTEMPT

Given an annotated program  $\langle P, L \rangle$ , we need automatic methods for: (1) computing the LC score of a given test suite  $TS$ , and (2) deriving a test suite achieving high LC-coverage. We propose first a black-box approach, reusing standard automatic testing tools through a *direct instrumentation* of  $P$ . This technique underlies previous works aiming at extending DSE coverage abilities [9], [15], [18], [20], [34], [35], [43]. While it allows for cheap LC score computation, it is far from efficient for ATG, mainly because of an exponential blow-up of the path space of the program.

##### A. Direct instrumentation

The *direct instrumentation*  $P'$  for  $\langle P, L \rangle$  consists in inserting for each label  $l \triangleq (loc, \varphi) \in L$  a new branching instruction  $I: if (\varphi) \{ \}$ ; such that all instructions leading to  $loc$  in  $P$  lead to  $I$  in  $P'$ , and  $I$  leads to  $loc$ . The transformation is depicted in Figure 4. When different labels

are attached to the same location, the new instructions are chained together in a sequence ultimately leading to *loc*.

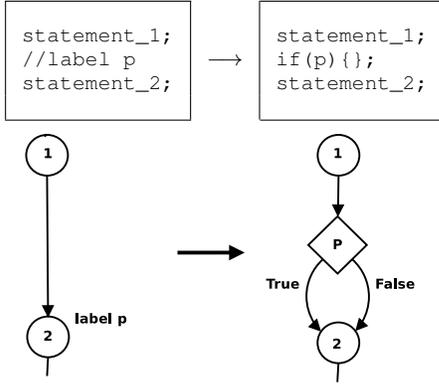


Fig. 4. Direct instrumentation  $P'$

Let us denote by **NTD** the set of test objectives over  $P'$  requiring to cover all New **T**hen-**D**ecisions introduced by the instrumentation. Direct instrumentation is obviously sound w.r.t. **LC** in the following sense.

*Theorem 3 (Soundness):* Given an annotated program  $\langle P, L \rangle$ , its direct instrumentation  $P'$  and a test suite  $TS$ , we have:  $TS \rightsquigarrow_{\langle P, L \rangle} \mathbf{LC}$  iff  $TS \rightsquigarrow_{P'} \mathbf{NTD}$ .

This is interesting for both **LC** score computation and ATG. Indeed, any ATG tool run on  $P'$  will produce a test suite  $TS$  covering **LC** for  $\langle P, L \rangle$  as soon as  $TS$  covers all branches of interest in  $P'$ . Besides, a slightly modified version of direct instrumentation, updating coverage information in the new then-branches, allows efficient coverage score computation.

*Theorem 4:* Given an annotated program  $\langle P, L \rangle$ , its direct instrumentation  $P'$  and a test suite  $TS$ , then the **LC** score of  $TS$  can be computed in time bounded by  $|TS| \cdot \text{maxtime}(\{P'(t) \mid t \in TS\})$ .

Computing **LC** score can be done independently from  $|L|$ . Regarding coverage score computation, **LC** is much closer to **DC** (each test  $t$  is executed only once) than it is to **WM** (each test  $t$  is executed *once per mutant*). While efficient score computation is a difficult issue in mutation testing, Theorems 2 and 4 show that the side-effect free subset of **WM** does enjoy efficient score computation.

### B. Drawbacks

So far, the direct instrumentation seems to perfectly suit our needs. Unfortunately, it is significantly inefficient for ATG. There are two main reasons for that.

- (†)  $P'$  is too complex: it exhibits many more behaviours than  $P$ , most of them being unduly complex for covering the labels we are targeting.
- (‡) DSE will naturally produce a test suite covering several times the same labels, which is useless since each label needs to be covered only once.

We formalize the first point (†) hereafter. We consider two dimensions in which  $P'$  is “too complex”: the size of the search space, denoted  $|Paths(P')|$ , and the shape of paths in

$Paths(P')$ . Let us call *label constraints* all additional branches  $\varphi$  and  $\neg\varphi$  introduced in  $P'$  compared to  $P$ , and let us denote by  $m$  the maximal number of labels *per* location in  $P$ . A single path  $\sigma \in P$  may correspond to up to  $2^{m \cdot |\sigma|}$  paths in  $P'$ , since each label of  $P$  creates a branching in  $P'$  and at most  $m$  such branchings can be found at each step of  $\sigma$ . Note also that the paths  $\sigma' \in P'$  corresponding to  $\sigma \in P$  have length bounded by  $m \cdot |\sigma|$ . Therefore they can pass through up to  $m \cdot |\sigma|$  label constraints, while (by definition)  $\sigma$  does not pass through any label constraint. Theorem 5 summarises these results.

*Theorem 5 (Non-tightness):* Given an annotated program  $\langle P, L \rangle$  and its direct instrumentation  $P'$ , let us assume that  $Paths(P)$  is bounded, that  $k$  represents the maximal length of paths in  $Paths(P)$  and that  $m$  is the maximal number of labels *per* location in  $P$ . Then:

- $|Paths(P')|$  can be exponentially larger than  $|Paths(P)|$  by a factor  $2^{m \cdot k}$ ;
- any  $\sigma' \in Paths(P')$  may carry up to  $m \cdot k$  (positive or negative) label constraints.

Both aspects are problematic for symbolic exploration of the search space: more paths means either more requests to a theorem prover (DSE) or a larger formula (BMC), while more constrained paths means more expensive requests.

## V. EFFICIENT ATG FOR LC

We describe in this section two main ingredients in order to obtain efficient ATG for **LC**: (1) a *tight instrumentation* avoiding all drawbacks of the direct instrumentation, and (2) a strong coupling of label coverage and DSE through *iterative label deletion*.

### A. Tight instrumentation

Given a label  $l \triangleq (loc, \varphi)$ , the key insights behind the tight instrumentation are the following:

- label constraint  $\varphi$  is useful only for covering  $l$ , and should not be propagated beyond that point;
- label constraint  $\neg\varphi$  is pointless w.r.t. covering  $l$ , and should not be enforced in any way.

Keeping these lines in mind, the instrumentation works as depicted in Figure 5: for each label  $(loc, \varphi)$ , we introduce a new instruction `if (nondet) {assert( $\varphi$ ); exit};` where `assert( $\varphi$ )` requires  $\varphi$  to be verified, `exit` forces the execution to stop and `nondet` is a non-deterministic choice<sup>1</sup>.

In the resulting instrumented program  $P^*$  (Figure 5, right column), when an execution reaches *loc*, it gives rise to two execution paths: the first one tries to cover the label by asserting  $\varphi$  and *stops right there*, the second one simply follows its execution *as it would do in P*, neither  $\varphi$  nor  $\neg\varphi$  being enforced.

Let us denote by **NA** the test objective over  $P^*$  requiring to cover all New **A**ssert introduced by the instrumentation (with condition evaluating to true). Tight instrumentation is sound w.r.t. **LC** in the following sense.

<sup>1</sup>Note that any DSE engine can simulate non-deterministic choices by an additional input array of (symbolic) boolean values.

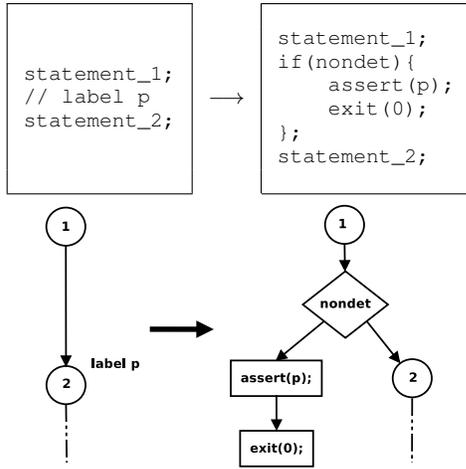


Fig. 5. Tight instrumentation  $P^*$

**Theorem 6 (Soundness):** Given an annotated program  $\langle P, L \rangle$ , its tight instrumentation  $P^*$  and a test suite  $TS$ , we have:  $TS \rightsquigarrow_{\langle P, L \rangle} \mathbf{LC}$  iff  $TS \rightsquigarrow_{P^*} \mathbf{NA}$ .

Interestingly, tight instrumentation does not show any of the issues reported in Theorem 5. The underlying reasons have been sketched at the beginning of Section V-A and are depicted in Figure 6. A single execution path in  $P$  going through  $n$  labels can give birth up to  $2^n$  paths in  $P'$  (left column), while it can create only  $n + 1$  paths in  $P^*$  (right column). Moreover, each path in  $P^*$  can go through at most one single positive label constraint, while a path  $\sigma'$  in  $P'$  can carry up to  $|\sigma'|$  (positive or negative) label constraints. These results are summarized in Theorem 7.

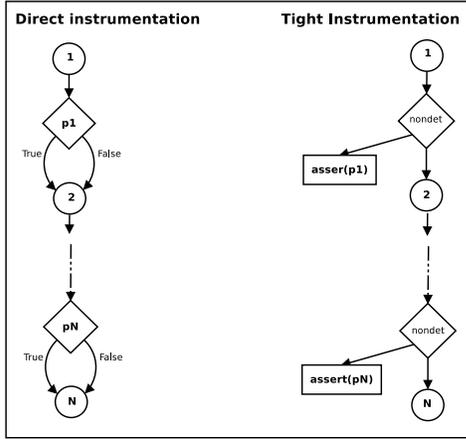


Fig. 6. Direct vs. tight instrumentation

**Theorem 7 (Tightness):** Given an annotated program  $\langle P, L \rangle$  and its instrumented version  $P^*$ , let us assume that  $Paths(P)$  is bounded, that  $k$  is the maximal length of paths in  $Paths(P)$  and that  $m$  is the maximal number of labels *per* location in  $P$ . Then  $P^*$  is tight in the following sense:

- $|Paths(P^*)|$  is linear in  $|Paths(P)|$  and  $m \cdot k$ ;
- any  $\sigma \in Paths(P^*)$  carries at most one label constraint.

*Proof:* The main reasons behind this result directly follow from the tight instrumentation, and have already been exposed just before Theorem 7. We can be more precise:  $|Paths(P^*)|$  is bounded by  $(m \cdot k + 1) \cdot |Paths(P)|$ . ■

Theorems 5 and 7 imply that any path-based program analysis conducted over  $P^*$  will have a much easier task than if conducted over  $P'$ , since  $P^*$  contains exponentially less paths and those paths are simpler.

### B. Iterative label deletion

We focus now on issue (‡) pointed out in Section IV-B. A DSE procedure launched on  $P^*$  tries to cover all paths from  $P^*$ , while we are only interested in covering branches corresponding to labels. Especially, it may try to cover path prefixes ending in an already-covered `assert( $\varphi$ )`. Whether they fail or not, these computations are redundant since Theorem 6 only requires that each new `assert` be covered once.

*Iterative label deletion (IDL)* consists in (conceptually) erasing a label constraint as soon as it is covered, so that it will not affect the subsequent path search. IDL requires to modify SE/DSE in the following way: each label  $l$  is equipped with a boolean variable  $b_l$  set to true iff  $l$  has already been covered during path exploration, and attempts to symbolically execute paths leading to  $l$  continue as long as  $b_l$  is false.

We present DSE with IDL over annotated programs in Algorithm 2, denoted  $DSE^*(\langle P, L \rangle)$ , where modifications w.r.t. standard SE/DSE are pointed out by  $(\star)$  marks. We assume that  $Paths(\langle P, L \rangle)$  is constructed in the following way: at each step, a run encountering a label  $l \triangleq (loc, \varphi)$  can either choose to go through  $l$  (enforcing  $\varphi$ ) and continue, or bypass  $l$  (no constraint) and continue.

---

#### Algorithm 2: Symbolic Execution with IDL

---

**Input:** an annotated program  $\langle P, L \rangle$  with finite set of paths  $Paths(P)$   
**Output:**  $TS$ , a set of pairs  $(t, \sigma)$  such that  $P(t) \rightsquigarrow_P \sigma$

- 1  $TS := \emptyset$ ;
- 2  $S_{paths} := Paths(\langle P, L \rangle)$ ;
- 3 **while**  $S_{paths} \neq \emptyset$  **do**
- 4     choose  $\sigma \in S_{paths}$ ;  $S_{paths} := S_{paths} \setminus \{\sigma\}$ ;
- 5     compute  $\phi_\sigma$ ;
- 6     **switch**  $solve(\phi_\sigma)$  **do**
- 7         **case**  $sat(t)$ :
- 8              $TS := TS \cup \{(t, \sigma)\}$ ;
- 9              $(\star)$  for all  $l$  covered by  $\sigma$ , do:  $b_l := 1$ ;
- 10             $(\star)$  remove from  $S_{paths}$  all  $\sigma'$  going through a label  $l$  s.t.  $b_l = 1$ ;
- 11         **case**  $unsat$ : skip;
- 12     **endsw**
- 13 **endsw**
- 14 **end**
- 15 **return**  $TS$ ;

---

For integration in a realistic SE/DSE setting with dynamic exploration of path space, we distinguish two flavors of IDL:

- IDL-1 a label is marked as covered only when it belongs to a path prefix being successfully solved. This is a purely symbolic approach.

IDL-2 a label is also marked when it is covered by a concrete execution, taking advantage of dynamic runs to delete several labels at once.

**Combining IDL with tight instrumentation.** Both variants of IDL can be combined with tight instrumentation either in a dedicated manner or in a black-box setting. Since dedicated implementations are straightforward, we focus hereafter on black-box implementations. An instrumentation enforcing IDL-1 over  $P^*$  is depicted in Figure 7. We add extra boolean variables for coverage, denoted  $b\_l$  where  $l$  is a label identifier. Yet, it is mandatory that the coverage information be global to the whole search rather than bound to a single path, and that the  $b\_l$  be treated concretely by DSE. For simplicity, we assume that it is achieved by putting the coverage information in an external file, accessed and modified through (concretized only) operations `read(b_l)` and `set_covered(b_l)`.

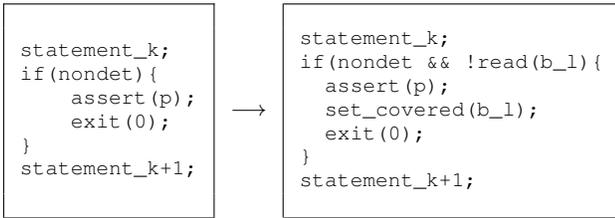


Fig. 7. IDL-1 variant of tight instrumentation  $P^*$

Enforcing IDL-2 in a black-box setting requires a fine-grained control over the DSE procedure. We need to be able to query the DSE engine for the next generated test datum. The procedure reuses the IDL-1-variant of tight instrumentation  $P^*$ , but each new generated test datum  $t$  is also run on the *direct instrumentation*  $P'$ . Unlike  $P^*$ ,  $P'$  does not exit at the first covered label, so all labels covered by  $t$  will be marked in the coverage file before the next test datum is searched for. The technique is depicted in Figure 8.

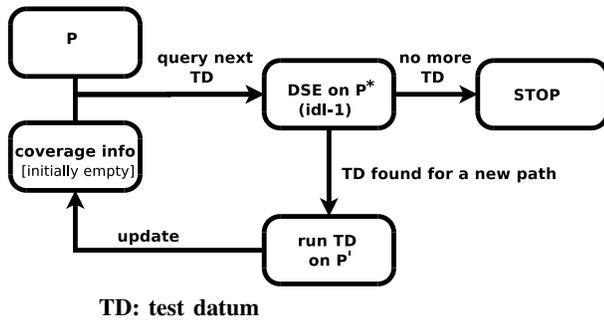


Fig. 8. IDL-2 variant for DSE

We denote by  $DSE^*(P^*)$  the combination of tight instrumentation and IDL presented in Figure 8, and we denote by  $DSE(\langle P, L \rangle)$  the variant of Algorithm 2 where lines 9-11 are removed. Considering only deterministic DSE techniques, the following result holds.

*Theorem 8 (Relative completeness):* Given an annotated program  $\langle P, L \rangle$  and its tight instrumentation  $P^*$ , then both  $DSE^*(\langle P, L \rangle)$  and  $DSE^*(P^*)$  cover as many labels as  $DSE(\langle P, L \rangle)$  does.

*Proof:* We first show that discarding paths going through an already covered label in  $DSE^*(\langle P, L \rangle)$  (cf lines 9-11 of

Algorithm 2) does not decrease label coverage w.r.t standard DSE on  $\langle P, L \rangle$ . Given a label  $l$ , then for any path  $\sigma \in Paths(\langle P, L \rangle)$  covering  $l$ , there exists by definition of  $Paths(\langle P, L \rangle)$  another path  $\sigma'$  passing through  $l$  and the same statements in  $P$  as  $\sigma$  while bypassing any other label, so that  $\sigma$  covers  $l$  iff  $\sigma'$  covers  $l$ . Since  $\sigma'$  cannot be discarded in  $DSE^*(\langle P, L \rangle)$  unless  $l$  is covered, it follows that the coverage cannot decrease, proving that  $DSE^*(\langle P, L \rangle)$  and  $DSE(\langle P, L \rangle)$  achieve the same coverage. Regarding coverage equivalence between  $DSE^*(P^*)$  and  $DSE(\langle P, L \rangle)$ , we reduce it to proving coverage equivalence between  $DSE^*(P^*)$  and  $DSE(P^*)$  through Theorem 6. We can then conclude with a similar proof technique. ■

## VI. IMPLEMENTATION & EXPERIMENTS

### A. Implementation

We have implemented tight instrumentation and iterative label deletion inside PATHCRAWLER [40]. The tool targets safety-critical C programs, with a strong focus on relative completeness guarantees. For example, the underlying constraint solver deals precisely with modular arithmetic, bitwise operations, floats and multi-level pointer dereferences. The DSE engine relies on a basic depth-first search and is highly optimised for programs with many infeasible paths [3].

Our implementation follows the description of Section V. We implement tight instrumentation and IDL-2 in a dedicated manner rather than a full black-box approach because PATHCRAWLER does not offer yet the required API for IDL-2. The search heuristics is mostly depth-first, but labels are handled as soon as possible.

### B. Experiments

Preliminary experiments have been conducted in order to investigate the following properties: (i) the relative gain of our two optimisations w.r.t. direct instrumentation, (ii) the overhead of leveraging DSE to LC. Evaluating the practical feasibility of label-based DSE over very large programs or its bug-finding power are left as future work<sup>2</sup>.

**Protocol.** We consider standard benchmark programs<sup>3</sup> taken from related works [9], [34], [32], mainly coming from the Siemens test suite (Tcas), the Verisec benchmark (get\_tag and full\_bad from Apache source code) and MediaBench (gd from libgd). We also consider three classes of labels simulating standard coverage criteria of increasing difficulty: CC, MCC and WM (cf. Section III-B). For WM, we mimic mutations introduced by MuJava [22] for operators AOIU, AOR, COR and ROR [2]<sup>4</sup>. Annotation is done manually<sup>5</sup>. Uncoverable labels (typically coming from equivalent mutants) are not discarded.

We compare the following algorithms:  $DSE(P)$  denotes standard DSE (witness),  $DSE(P')$  denotes standard DSE on direct instrumentation,  $DSE(P^*)$  denotes standard DSE on tight instrumentation and  $DSE^*(P^*)$  denotes DSE with iterative label deletion run on tight instrumentation. The DSE engine runs in deterministic mode, generating the same concrete

<sup>2</sup>Bug-finding power of criteria in Section III is already extensively studied.

<sup>3</sup>Available at <http://sebastien.bardin.free.fr/benchs-icst.zip>

<sup>4</sup>These operators are considered very powerful in practice [27], [39].

<sup>5</sup>The cost of automatic annotation is negligible w.r.t. the cost of DSE.

values from one run to the other. Time-out for solver is set to 1 min, time-out for test generation is set to 1h30. Experiments are performed on an Intel Core2 Duo 2.40GHz, 4GB of RAM.

We record the following information: number of paths explored by the search, computation time and achieved coverage. The number of paths is a good measure for comparing the complexity of the different search spaces, and therefore to assess both the “cost” of leveraging DSE to labels and the benefits of our optimisations. Coverage score together with computation time indicate how practical label-based DSE is.

It must be highlighted that PATHCRAWLER does not stop until all feasible paths are explored. This strategy gives us a good estimation of the size of the path space, however in practice it would be wiser to implement a label-based stopping criterion. Hence, from a feasibility point of view, results reported here are too pessimistic.

**Results.** A representative subset of results is presented in Table I. First, note that when no time-out occurs, direct instrumentation and both variants of tight instrumentation achieve the same coverage, and that this coverage is high (>90% on 17/25 examples). We also observe that direct instrumentation yields a significant overhead, confirming previous work [18]:  $DSE(P')$  has four time-outs (TO) while  $DSE(P)$  has none, time-overhead goes up to 122x (excluding TO), growth of the path-space reaches 50x.

On the other hand, tight instrumentation  $DSE^*(P^*)$  yields only a very reasonable overhead w.r.t. standard DSE: no time-out is reported, time-overhead is kept under 7x with an average of 2.4x, growth of the path-space is limited to 3x. On some examples, tight instrumentation performs remarkably better than direct instrumentation (94s vs TO on `gd_5-wm`). Interestingly,  $DSE^*(P^*)$  does perform better than standard DSE (up to 2x) on a few examples with very few additional paths. We conjecture that additional label constraints may sometimes greatly simplify the solving process, but it must be investigated further.

Finally, as expected,  $DSE(P^*)$  stands between  $DSE(P')$  and  $DSE^*(P^*)$  for the number of paths. Results are more mitigated for computation time, where  $DSE(P^*)$  is slower than  $DSE(P')$  on several examples, probably due to the fact that our implementations of  $DSE(P^*)$  and  $DSE^*(P^*)$  are not optimal.

**Conclusion.** These experiments confirm our formal predictions:

- fully-optimised DSE performs significantly better on difficult programs than the direct instrumentation, both in terms of search space and computation time;
- the overhead w.r.t. standard DSE turns out to be always acceptable, and often very low.

These results suggest that DSE can be efficiently leveraged to **LC** coverage thanks to our optimisations. Further experiments on larger programs are required to fully confirm that point.

## VII. BEYOND TEST DATA GENERATION

Section III proves that **LC** is a powerful coverage criterion, encompassing many standard criteria and a large subset of weak mutations. Section V and Section VI demonstrate the feasibility of efficient ATG for **LC**, with a cost-effective integration in DSE. We also sketched in Section IV how to perform

			DSE(P) (witness)	DSE(P')	DSE(P*)	DSE*(P*)
trityp 50 loc	cc	#paths time cover	35 1.3s	183 1.6s 24/24	83 2s 24/24	46 4.5s 24/24
	mcc	#paths time cover	35 1.3s	337 1.9s 28/28	110 3s 28/28	66 2.1s 28/28
	wm	#paths time cover	35 1.3s	x x x	506 12s 120/129	48 5.1s 120/129
4balls 35 loc	wm	#paths time cover	7 1.2s	<b>195</b> 1.9s 56/67	75 2.1s 56/67	<b>23</b> 2.1s 56/67
utf8-3 108 loc	wm	#paths time cover	134 1.4s	1,379 4.2s 55/84	626 4.3s 55/84	313 3.8s 55/84
utf8-5 108 loc	wm	#paths time cover	680 2s	<b>11,111</b> 40s 82/84	3,239 24s 82/84	<b>743</b> 8.1s 82/84
utf8-7 108 loc	wm	#paths time cover	3,069 5.8s	<b>81,133</b> <b>576s</b> 82/84	14,676 110s 82/84	<b>3,265</b> <b>35s</b> 82/84
tcas 124 loc	cc	#paths time cover	2,787 2.9s	3,508 3.6s 10/10	3,508 5s 10/10	2,815 3.4s 10/10
	mcc	#paths time cover	2,787 2.9s	3,988 4.2s 11/12	3,988 5.2s 11/12	3,059 3.9s 11/12
tcas' 124 loc	wm	#paths time cover	4,420 5.6s	<b>300,213</b> <b>662s</b> 101/111	20,312 120s 101/111	<b>6,014</b> <b>27s</b> 101/111
replace 100 loc	wm	#paths time cover	866 2s	<b>87,498</b> <b>245s</b> 70/79	6,420 64s 70/79	<b>2,347</b> <b>14s</b> 70/79
full_bad 219 loc	cc	#paths time cover	2,563 5s	5,148 8s 12/16	5,129 14s 12/16	3,209 7s 12/16
	mcc	#paths time cover	2,563 5s	12,360 19s 24/39	12,296 32s 24/39	7,043 19s 24/39
	wm	#paths time cover	2,593 5s	19,336 35s 34/46	10,610 40s 34/46	5,414 19s 34/46
get_tag-5 240 loc	cc	#paths time cover	11,833 60s	40,102 210s 20/20	22,669 651s 20/20	11,843 64s 20/20
	mcc	#paths time cover	11,833 60s	41,605 100s 26/26	23,794 510s 26/26	11,848 48s 26/26
	wm	#paths time cover	11,833 61s	58,646 140s 44/47	28,919 719s 44/47	11,856 51s 44/47
get_tag-6 240 loc	cc	#paths time cover	76,456 3,011s	<b>TO</b>	<b>TO</b>	76,468 <b>1,512s</b> 20/20
	wm	#paths time cover	76,456 3,011s	<b>TO</b>	<b>TO</b>	76,481 <b>1,463s</b> 44/47
gd-5 319 loc	cc	#paths time cover	14,516 52s	18,220 66s 36/36	17,018 91s 36/36	14,605 59s 36/36
	mcc	#paths time cover	14,516 51s	20,261 71s 29/36	18,799 101s 29/36	15,201 80s 29/36
	wm	#paths time cover	14,516 50s	<b>TO</b>	<b>TO</b>	14,607 <b>94s</b> 62/63
gd-6 319 loc	cc	#paths time cover	107,410 3,740s	131,726 3,816s 36/36	125,024 5,534s 36/36	107,500 2,945s 36/36
	mcc	#paths time cover	107,410 3,740s	144,840 3,822s 29/36	137,328 6,281s 29/36	111,208 3,447s 29/36
	wm	#paths time cover	107,410 3,740s	<b>TO</b>	<b>TO</b>	107,521 <b>2,232s</b> 63/63

TO: time-out (5,400 sec) x: crash due to a bug in the underlying solver

TABLE I. EXPERIMENTAL RESULTS FOR ATG

cheap **LC** score computation. Everything put together, labels form the basis of a very powerful framework for automatic testing, handling many different criteria in a uniform fashion. We describe such a view in Figure 9. Starting from a program  $P$  and a testing criterion  $C$ , a predefined labelling function  $\psi_C$  creates the  $C$ -equivalent annotated program  $\langle P, L \rangle$  (Theorems 1 and 2). Then, we can perform efficient **LC** score computation and **LC**-based ATG through instrumentation (Theorems 4 and 6). Finally, static analysis techniques can be used on  $P^*$  in order to detect uncoverable labels, i.e. labels  $l \triangleq (loc, \varphi)$  for which there is no test datum  $t$  such that  $t \rightsquigarrow_{\psi_C(P)} l$ . Static detection of uncoverable labels can help ATG tools by avoiding wasting time on infeasible objectives, as was observed in the case of mutation testing [19].

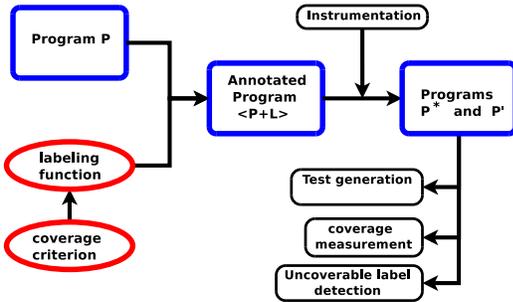


Fig. 9. **LC**-coverage framework

## VIII. RELATED WORK

**Leveraging DSE to higher coverage criteria.** The need for enhancing DSE with better coverage criteria has already been pointed out in active testing (a.k.a assertion-based testing) [9], [15], [20], Mutation DSE [32], [33] and Augmented DSE [18], [35], [43]. The present work generalizes these results and proposes ways of taming the potential blow-up, resulting in an effective support of advanced coverage criteria in DSE with only a small overhead.

Active testing targets run-time errors by adding explicit branches into the program. It is similar to the Run-Time Error Coverage criterion presented in Section III. Labels are a more general approach. The direct instrumentation  $P'$  for this criterion is mostly equivalent to  $P^*$  since additional branches can only trigger errors and stop the execution. Yet, active testing could benefit from the IDL optimisation. In that case only the IDL-1 flavour makes sense since an execution cannot cover two different run-time errors. Finally, since most test objectives are (hopefully!) uncoverable for Run-Time Error Coverage, some approaches aim at combining DSE with static detection of uncoverable targets [9]. They can be reused for labels, and should be useful when many labels are uncoverable.

Following Offut *et al.* [12], Papadakis *et al.* show that **WM** can be reduced to branch coverage through the use of a variant of Mutant Schemata [38]. This is pretty similar to the direct encoding  $P'$  mentioned here. They propose essentially two variations of DSE for mutation testing: a black-box approach [32] based on a direct encoding similar to our DSE( $P'$ ) scheme, and a more ad hoc approach [32] preventing reuse of existing DSE tools but offering several optimisations.

Papadakis *et al.* propose a variant of IDL, a dedicated search heuristic based on shortest paths [31] and an improvement of the direct encoding through the use of mutant identifiers (following exactly Mutant Schemata). On the one hand, it ensures that a given path cannot go through several *different* mutants, on the other hand there is still an exponential blow-up of the search space in the worst case, and IDL cannot cover more than one mutant at once.

Augmented DSE [18] is a variant of direct instrumentation. Several coverage criteria are encoded, getting results similar to those of Section III-B, yet the side-effect free subset of **WM** is not identified. Experiments [18, Table 2] report an average time-overhead of 272x, going up to >2,000x. That confirms the strong benefits of our optimisations, that yield a maximal overhead of 7x.

We give a more generic view of the problem, identifying labels and annotated programs as the key concept underlying the approach. We also clearly identify the limits and hypotheses of the method by defining the side-effect free fragment of **WM**, proving soundness of direct instrumentation and providing a formalization of the path space “complexification” induced by direct instrumentation. Most important, we propose the tight instrumentation which completely prevents complexification. Finally, our optimisations can be implemented in a pure black-box setting and we do not impose anything on the search heuristics, keeping room for future improvements.

**Labels and optimized DSE.** The label-specific optimisations described here can be freely mixed with other DSE optimisations. It is left as future work to explore which optimisations turn out to be the most effective for labels. As already stated, combining static discovery of uncoverable labels with DSE [9] could be useful for often-uncoverable labels, such as those generated for Run-Time Error Coverage or **MCC**. Another promising direction is to adapt DSE search heuristics [41] by taking advantage of the dissimilarities between labels and branches, possibly getting inspiration from [31].

The IDL optimisation shows some similarities with Look-Ahead pruning (LA) [4], [6]. Basically, LA takes advantage of (global) static analysis to prune path prefixes which cannot reach any uncovered branches. On  $P^*$ , IDL-1 is a very specific (but cheap) case of LA while IDL-2 is orthogonal: LA prunes those “label paths” pruned by IDL-1 plus other paths leading only to already covered labels, while IDL-2 prunes several “label paths” at once thanks to dynamic analysis.

**Automation of mutation testing.** Mutation coverage [11], [28] has been established as a powerful criterion through several experimental studies [1], [28]. Yet, it is very difficult to automatize. Even mutation score computation is expensive in practice if not done wisely. Weak mutations [17] relax mutation coverage by abandoning the “propagation step”, making **WM** easier to compare with standard criteria and easier to test for. **WM** has been experimentally proved to be almost equivalent to strong mutations [26], and from a theoretical point of view **WM** subsumes many other criteria [29].

The few existing symbolic methods for mutation-based ATG are based on the encoding proposed by Offutt *et al.* and have already been discussed [12], [34], [33]. The Mutation Schemata technique [38] was originally developed in order to

factorize the compilation costs of hundreds of similar mutants. Static analysis has been proposed for the “equivalent mutant detection” problem [24], [23] in a way similar to what is sketched in Section VII.

The side-effect free fragment of **WM** presented in this paper seems to be a sweet spot of mutation testing: it is amenable to efficient automation and still very expressive. It is left as future work to identify if something essential is lost within this fragment. Finally, our encoding of **WM** into **LC** is orthogonal to and can be combined with some of the many techniques developed for efficient mutation testing, such as operator reduction [27], [39] or smart use of operators [19].

## IX. CONCLUSION

Label coverage is a new testing criterion which appears to be both expressive and amenable to efficient automation. Some of the ideas behind labels underly previous work by other teams. We generalise them, propose ways of taming the potential complexification of the path space and provide both formal and experimental evidence. Especially, we have shown how to extend DSE for label coverage in a black-box manner with only a slight overhead. Experiments show that our optimisations yield very significant improvements.

This work bridges part of the gap between symbolic ATG techniques and coverage criteria. On the one hand, we show that DSE techniques can be cheaply extended to support more advanced testing criteria, including side-effect free weak mutations. On the other hand, we identify a powerful criterion amenable to efficient automation, both in terms of ATG and coverage score computation.

## REFERENCES

- [1] J. H. Andrews, L. C. Briand, Y. Labiche: Is mutation an appropriate tool for testing experiments? In: ICSE 2005. IEEE
- [2] P. Ammann, A. J. Offutt: Introduction to software testing. Cambridge University Press, New York (2008)
- [3] B. Botella, M. Delahaye, S. Hong-Tuan-Ha, N. Kosmatov, P. Mouy, M. Roger, N. Williams: Automating Structural Testing of C Programs: Experience with PathCrawler. In: AST 2009. IEEE
- [4] S. Bardin and P. Herrmann. Pruning the search space in path-based test generation. In: ICST 2009. IEEE
- [5] S. Bardin, P. Herrmann. OSMOSE: Automatic Structural Testing of Executables. *Softw. Test., Verif. Reliab.* 21(1): 29-54(2011)
- [6] J. Burnim, K. Sen. Heuristics for Scalable Dynamic Test Generation. In: ASE 2008. IEEE
- [7] C. Cadar, D. Dunbar, D. Engler: KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs. In: OSDI 2008. Usenix Association (2008)
- [8] C. Cadar, V. Ganesh, P. M. Pawlowski, D. L. Dill and D. R. Engler. EXE: automatically generating inputs of death. In: CCS 2006. ACM
- [9] O. Chebaro, N. Kosmatov, A. Giorgetti, J. Julliard: Program slicing enhances a verification technique combining static and dynamic analysis. In: SAC 2012. ACM
- [10] E. M. Clarke, D. Kroening, F. Lerda: A Tool for Checking ANSI-C Programs. In: TACAS 2004. Springer
- [11] R. A. DeMillo, R. J. Lipton, A. J. Perlis: Hints on test data selection: Help for the Practicing Programmer. *Computer*, 11(4), 34-41
- [12] R. A. DeMillo, A. J. Offutt: Constraint-Based Automatic Test Data Generation. *IEEE Trans. Software Eng.* 17(9), 1991
- [13] P. Godefroid, N. Klarlund and K. Sen. DART: Directed Automated Random Testing. In: PLDI 2005. ACM
- [14] P. Godefroid, M. Y. Levin and D. Molnar. Automated Whitebox Fuzz Testing. In: NDSS 2008.
- [15] P. Godefroid, M. Y. Levin, D. Molnar: Active property checking. In: EMSOFT 2008. ACM
- [16] P. Godefroid, M. Y. Levin, D. A. Molnar: SAGE: whitebox fuzzing for security testing. *Commun. ACM* 55(3): 40-44 (2012)
- [17] W. E. Howden: Weak mutation testing and completeness of test sets. In: *IEEE Transactions on Software Engineering*, 8(4). 1982
- [18] K. Jamrozik, G. Fraser, N. Tillmann and J. de Halleux. Generating Test Suites with Augmented Dynamic Symbolic Execution. In: TAP 2013.
- [19] R. Just, G. M. Kapfhammer and F. Schweiggert. Do Redundant Mutants Affect the Effectiveness and Efficiency of Mutation Analysis? In: ICST 2012. IEEE
- [20] B. Korel, A. M. Al-Yami: Assertion-Oriented Automated Test Data Generation. In: ICSE 1996. IEEE
- [21] J. C. King. Symbolic execution and program testing. *Communications of the ACM*, 19(7), July 1976.
- [22] Y. S. Ma, A. J. Offutt, Y. R. Kwon: MuJava: a mutation system for java. In: ICSE 2006. ACM
- [23] S. Nica, F. Wotawa: Using Constraints for Equivalent Mutant Detection. In: workshop Formal Methods in the Development of Software (2012)
- [24] A. J. Offutt, W. M. Craft: Using Compiler Optimization Techniques to Detect Equivalent Mutants. *Softw. Test., Verif. Reliab.* 4(3). 1994
- [25] A. J. Offutt: Investigations of the Software Testing Coupling Effect. *ACM Trans. Softw. Eng. Methodol.* 1(1): 5-20 (1992)
- [26] A. J. Offutt, S. D. Lee: An Empirical Evaluation of Weak Mutation. *IEEE Trans. Software Eng.* 20(5): 337-344 (1994)
- [27] A. J. Offut, G. Rothermel, C. Zapf: An experimental evaluation of selective mutation. In: ICSE 1993. IEEE
- [28] A. J. Offutt, R. H. Untch: Mutation 2000: uniting the orthogonal. In: Mutation testing for the new century. Kluwer Academic Publisher, 2001
- [29] A. J. Offutt, J. Voas: Subsumption of Condition Coverage Techniques by Mutation Testing. Tech. Report ISSE-TR-96-01, Dpt. Information and Software System Engineering, George Mason Univ., 1996
- [30] G. C. Necula and S. Mcpeak and S. P. Rahul and W. Weimer CIL: Intermediate language and tools for analysis and transformation of C programs. In: CC 2002. Springer
- [31] M. Papadakis, N. Malevris: An Effective Path Selection Strategy for Mutation Testing. In: APSEC 2009. IEEE
- [32] M. Papadakis, N. Malevris: Automatic Mutation Test Case Generation via Dynamic Symbolic Execution. In: ISSRE 2010. IEEE
- [33] M. Papadakis, N. Malevris: Automatically performing weak mutation with the aid of symbolic execution, concolic testing and search-based testing. *Software Quality Journal* 19(4), 2011
- [34] M. Papadakis, N. Malevris, M. Kallia: Towards Automating the Generation of Mutation Tests. In: AST 2010 (with ICSE 2010).
- [35] R. Pandita, T. Xie, N. Tillmann and J. de Halleux. Guided test generation for coverage criteria. In: ICSM 2010. IEEE
- [36] K. Sen, D. Marinov and G. Agha. CUTE: A Concolic Unit Testing Engine for C. In: ESEC/FSE 2005. ACM
- [37] N. Tillmann and J. de Halleux. Pex-White Box Test Generation for .NET. In: TAP 2008. Springer
- [38] R. H. Untch, A. J. Offutt, M. J. Harrold: Mutation Analysis Using Mutant Schemata. In: ISSA. ACM (1993)
- [39] W. E. Wong, A. P. Mathur: Reducing the cost of mutation testing: An empirical study. In: *Journal of Systems and Software*, 31(3), 185-196.
- [40] N. Williams, B. Marre and P. Mouy. On-the-Fly Generation of K-Path Tests for C Functions. In: ASE 2004. IEEE (2004)
- [41] T. Xie, N. Tillmann, P. de Halleux and W. Schulte. Fitness-Guided Path Exploration in Dynamic Symbolic Execution. In: DSN 2009. IEEE (2009)
- [42] H. Zhu, P. A. V. Hall and J. H. R. May. Software Unit Test Coverage and Adequacy. In: *ACM Computing Surveys*, vol. 29(4), 1997
- [43] L. Zhang, T. Xie, L. Zhang, N. Tillmann, J. de Halleux and H. Mei. Test generation via Dynamic Symbolic Execution for mutation testing. In: ICSM 2010. IEEE