

PROSECUTOR: Bayesian Counterfactual Fault Localization

ANONYMOUS AUTHOR(S)

Bayesian reasoning has emerged as an exciting approach to fault localization, where the introduction of errors and their subsequent propagation through faulty execution traces is treated as a stochastic process. One can then perform Bayesian inference on a probabilistic model encoding the program execution to associate individual statements and values with a posterior suspicion/belief of being erroneous. In this paper, we propose a new graph representation that effectively models error propagation through failing program executions. This structure, which we call the Error Propagation Graph (EPG), extends prior probabilistic approaches by incorporating richer interprocedural relationships and accounting for the influence of unexplored control-flow branches that may affect variable values. We also show how EPGs can be constructed efficiently and compactly, and how this structure enables the selection of a set of counterfactual experiments, each involving artificially flipping a suspicious branch predicate at runtime and observing its downstream effect on the test outcome. The results of these experiments provide additional evidence that can be incorporated back into the EPG to confirm or refute the model's suspiciousness estimates.

We have implemented this technique in a tool which we call PROSECUTOR, and evaluated it on 470 buggy versions of 13 projects from the Defects4J benchmark suite. Our experimental evaluation shows that 40% of the true fault locations are identified within just the three lines of code believed to be the most suspicious. The technique also turns out to be significantly more effective than a diverse set of baselines, and successfully identifies at least 12%, 10%, 15%, and 19% *more* true fault locations than each of the baselines within its predictions for the top 1, 3, 5, and 10 most suspicious program statements, respectively.

ACM Reference Format:

Anonymous Author(s). 2026. PROSECUTOR: Bayesian Counterfactual Fault Localization. 1, 1 (March 2026), 27 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

The fault localization problem may be stated as follows: Given a faulty program whose bug is witnessed by a failing test case, rank, in order of suspicion, the lines of code that are most likely to be the location of the bug. Starting with Jones and Harrold's groundbreaking work on Tarantula [13], this problem has been the subject of a rich body of research in software engineering.

Techniques range from the so-called *spectrum-based fault localization (SBFL)*—which monitor test coverage to identify parts of the program that are disproportionately executed by failing tests and infrequently executed by passing tests [6, 7, 13, 25, 32]—to *mutation-based fault localization (MBFL)*, which apply a series of tentative modifications to the program and observe their downstream effects on test outcomes [24, 26].

These approaches have complementary strengths and weaknesses: Spectrum-based methods are lightweight (they simply need some form of code coverage monitoring), but rely on access to a comprehensive suite of tests to be effective. Moreover, their reliance on aggregate statistics gives them limited discriminatory power. In contrast, mutation-based techniques provide actionable evidence for possible locations of the bug, but because of the large space of mutations and the need to repeatedly rerun the entire test suite, their scalability is limited in practice.

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

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2026/3-ART

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

50 Another promising approach to fault localization was recently introduced by Zeng et al. in
51 SmartFL [38], which constructs a probabilistic model of program execution using data and control
52 dependencies. SmartFL then assigns each statement a prior probability of being faulty and leverages
53 observations from test executions to perform marginal inference on the constructed graph and
54 calculate posterior probabilities that these individual statements are buggy.

55 Despite demonstrating the effectiveness of probabilistic approaches for fault localization, SmartFL
56 struggles to effectively localize faults in several scenarios. Specifically, the technique has difficulty
57 identifying faults when bugs are deep within the test execution and their effects propagate across
58 multiple procedures. Moreover, SmartFL's reliance solely on data and control dependencies is
59 insufficient to fully capture errors arising from incorrect evaluations of branch predicates that
60 prevent expected variable updates. These limitations highlight the challenges of effectively modeling
61 error propagation in programs with long execution chains and subtle control-flow dependencies.

62 With this background, our paper extends SmartFL's idea in two ways: *First*, we design a new
63 graph structure which we call it the error propagation graph (EPG) to more effectively model the
64 propagation of errors through failing executions. In addition to direct dependencies between the
65 site of a variable definition and the point of its use, we also, for example, include edges in the EPG
66 from unexplored branches that contain unexecuted variable assignments. Such careful modeling
67 is crucial, and as we will see in our experimental evaluation, our modeling contributes to a 15%
68 increase in the number of bugs localized by examining just 5 lines of code over SmartFL [35], the
69 previous state-of-the-art probabilistic baseline. This also contributes to a 4.2% median EXAM score
70 which is less than half of that achieved by SmartFL (10.2%).

71 *Next*, by attaching conditional probability distributions to each vertex in this graph, we view
72 the EPG as a Bayesian network. We perform Bayesian inference over the constructed network and
73 compute posterior suspiciousness scores for each individual branch predicate evaluated at runtime.
74 *Finally*, we use these posterior probabilities to identify a set of *counterfactual* experiments. Each
75 of these experiments involves artificially changing the value of the branch predicate at runtime
76 and observing its effect on eventual test outcome. Intuitively stated: If a branch condition was
77 incorrectly evaluated, then it is possible that artificially flipping its value might cause the failing
78 test to succeed. Alternatively, if a statement is buggy, then it is possible that rerouting control
79 flow around this statement prevents test failure. These experiments therefore help us to confirm
80 or refute our belief in the correctness/incorrectness of the branch conditions and their dependent
81 statements, providing an additional source of information that can be incorporated into the ranking
82 process.

83 Due to the availability of the Soot framework [30] and Just et al.'s Defects4J dataset of real-
84 world bugs [14], we implemented our approach to locate bugs in Java programs. In doing so, we
85 had to solve several technical challenges: First, marginal inference is famously intractable [17].
86 We therefore have to suitably compress the constructed EPG in order to use standard inference
87 algorithms such as loopy belief propagation [16]. Our solution involves an instrumentation-guided
88 loop identification algorithm to identify and compress loops in execution traces.

89 The second consideration involves determining which counterfactual experiments to conduct
90 and how to incorporate their results back into the probabilistic model. We restrict our analysis to
91 the 20 most suspicious branch conditions identified by the initial Bayesian inference run. In order
92 to leverage information from these experiments, we extend the initial probabilistic model with
93 virtual variables and edges. This allows us to model the impact of flipping the branch condition on
94 other dependent statements as well.

95 We call our implementation PROSECUTOR and evaluated it on faulty versions of 13 projects from
96 the Defects4J dataset [14]. We compare PROSECUTOR to 8 baselines: 4 from SBFL family, 2 from MBFL
97 family, and 2 state-of-the-art methods, SmartFL and DepGraph [27], which leverages probabilistic
98

models and graph neural networks (GNN), respectively. We observed that PROSECUTOR identifies *at least* 15% and 19% more true faulty statements within its top-5, and top-10 predictions than any of the baselines under consideration. Moreover, we observe that, unlike SBFL and MBFL, whose effectiveness strictly relies on the presence of a high-coverage test suite with both passing and failing tests, our approach can effectively localize faults without even consulting the passing traces.

In summary, this paper makes the following contributions:

- (1) A systematic approach to model runtime program dependencies in object-oriented languages using a data structure called the error propagation graph (EPG).
- (2) A new fault localization approach with a central probabilistic model of error introduction and propagation, which leverages counterfactual analysis to reason about the location of faults.
- (3) Optimizations, including loop compression and eliminating redundant dependencies, which enable scaling the technique to real-world software.
- (4) Empirical evaluation on buggy versions of 13 projects from the Defects4J benchmark to showcase that PROSECUTOR outperforms existing approaches in identifying faults.

2 Motivating Example

Consider the buggy Java program in Figure 1. As shown in the assertions on Lines 32 and 35, the goal of this function is to pad the left side of the input string until it reaches the desired length. In the case of the failing test, the incorrect assignment to the `padChars` variable leads to an unexpected out-of-bounds array access on Line 26. In contrast, the passing test case never uses the erroneous variable, and correctly returns with the expected output on Line 22.

We argue that neither spectrum-based nor mutation-based approaches satisfactorily localize this bug: In particular, as can be seen from the trace annotations, Lines 24–26 are the only lines of code that are uniquely executed by the failing test, but none of them are buggy. As such, spectrum-based approaches cannot distinguish between the remaining statements. On the other hand, MBFL techniques typically focus on mutating operators rather than operands (which require careful consideration of types and contexts to safely mutate) and are therefore unable to mutate the buggy Line 9. In addition, MBFL techniques such as Metallaxis also consider the effect of the mutation on test *behavior*, rather than simply the eventual *outcome*.

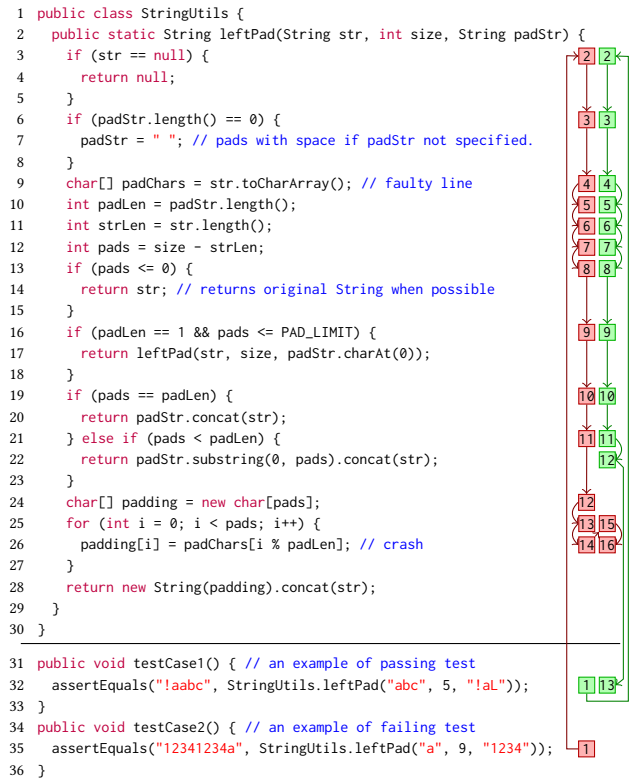


Fig. 1. Example implementation of the `leftPad` utility function, adapted from the Apache Commons Lang library [3]. The program contains a bug on Line 9, which should instead be `padChars = padStr.toCharArray()`. Graph annotations on the right side describe the execution flow in each test case.

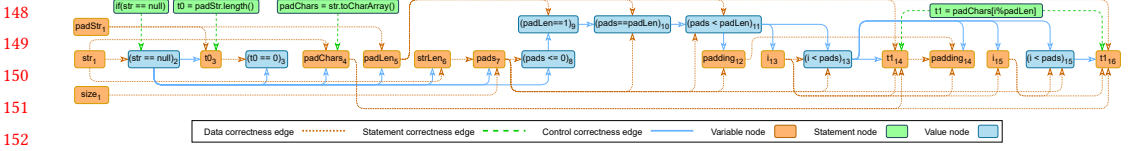


Fig. 2. Snippet of the error propagation graph (EPG) for the trace obtained from `testCase2` in Figure 1.

In this context, mutating the `if`-statements on Lines 3, 6, 13, 16, 19, and 21 may provide opportunities for premature return, thereby eliminating the crash on Line 26, and instead causing an assertion violation at Line 35. The algorithm therefore characterizes them as “successful” mutations, despite being irrelevant to the bug, leading to an unsatisfactory ranking of possible fault locations.

2.1 Postmortem of the Fault

Line 26 may be regarded as consisting of two operations: the first which loads and computes the value `padChars[i % padLen]`, and the second which stores this value in `padding[i]`. It is the first operation, i.e., the computation of `padChars[i % padLen]`, which results in the crash. By examining the source code and crash context, we conclude that one of the following *must* be the case: (a) the load operation on Line 26 is faulty, or (b) one of the operand variables, `padChars`, `i` and `padLen`, was incorrectly computed, or that (c) control must never have reached Line 26, forcing us to examine whether the loop guard on Line 25 was correctly evaluated in the last iteration of this loop. This backward analysis of program dependencies leads us to suspect and examine four lines of code: Lines 9, 10, 25, and 26. Much of our paper follows from the following more general proposition:

Consider the sequence of statements, t_1, t_2, \dots, t_n executed along a program trace. If executing a statement t_i introduces or propagates an error into the program state, then either: (a) t_i is itself buggy, or (b) one of its operands was incorrectly computed, or (c) the branch condition that controls the execution of t_i was incorrectly evaluated.

Our approach leverages program dependencies to construct a model of error propagation instead of just using coarse-grained coverage information. We encode the dependencies between the correctness / incorrectness of runtime values using a data structure that we call the *error propagation graph (EPG)*. A fragment of such a graph is illustrated in Figure 2.

At a high level, the error propagation graph contains three kinds of nodes: statement nodes, value nodes, and variable nodes respectively. Each statement node represents the correctness of the corresponding program statement, while variable nodes represent the correctness of the variables defined along the program trace. We also associate statements which do not define new variables but rather only evaluate expressions in order to change control flow, with so-called value nodes corresponding to the correctness of the expressions evaluated by these statements at runtime.

For example, the node `padChars4` in the graph of Figure 2 indicates that in the failing trace, the correctness of the value assigned to `padChars` in Step 4 depends on the correctness of the statement `padChars = str.toCharArray()` in Line 9 and the value of `str` defined at the method’s entry point. It also depends on whether the branch condition at Step 2 is evaluated correctly. Furthermore, the node illustrates that an error which reaches the execution Step 4 and taints the value of variable `padChars` may subsequently propagate and affect the value of the temporary variable `t1` in Steps 14 and 16. For the purpose of exposition and to avoid cluttering the graph, we only show a part of statement nodes which affect values such as `padChars4`, `t114`, and `t116`. Note however, that the complete EPG contains similar nodes for *all* other runtime values.

Another variable node of interest, $t1_{16}$, corresponds to the value temporarily created by the load operation at Step 16, *had it not crashed*. Following the earlier discussion above, the correctness of this variable depends on the correctness of the load statement, $padChars_4$, $padLen_5$, and i_{15} . It also depends on whether the branch condition ($i < pads$) at Step 15 was evaluated correctly.

Observe that the EPG effectively constrains the manner in which errors may be introduced or propagated along the execution trace. If a statement is incorrect, then forward analysis from the corresponding statement node will uncover all values and variables that might have been incorrectly computed. Conversely, if a value is found to be incorrect, then the erroneous statement must reside in the backward cone of influence from this node in the EPG.

Observe also that the subgraph corresponding to the value and variable nodes resembles the program dependence graph (PDG), which is constructed by integrating control and data dependencies. The key difference between conventional PDGs and this subgraph of the EPG is that the nodes in PDGs are program statements and the edges represent control and data dependencies between them. In contrast, the blue and yellow nodes in Figure 2 depict runtime values from each step of the execution trace and the edges represent correctness dependencies between them. Also, note that the EPG can be generalized to the setting of having multiple test executions. In these situations, the subgraphs corresponding to different test executions would share the same set of statement nodes, but would have their own separate sets of value and variable nodes, representing the runtime behavior of the corresponding execution.

2.2 The Probabilistic Model

When faced with buggy code, software engineers are initially unsure of the exact location of the fault, but might be suspicious of different program statements. Another key aspect of the error propagation graph is that it enables us to quantitatively reason about correlations between the correctness of various program elements.

For example, consider the node $padLen_5$. If all of its predecessors in the EPG were correctly evaluated, then this node must also be correctly computed. I.e.,

$$\Pr(padLen_5 \mid h_0) = 1,$$

where $h_0 = padStr_1 \wedge (str == null)_2 \wedge padLen = padStr.length()$ corresponds to the event that none of its predecessors were incorrect. On the other hand, this variable might also have been incorrectly computed for multiple reasons. If, for example, an incorrect input string $padStr$ was provided, then we declare:

$$\Pr(padLen_5 \mid h_1) = 0.85,$$

where $h_1 = \neg padStr_1 \wedge (str == null)_2 \wedge padLen = padStr.length()$. Notice that the statement has some degree of “*fault tolerance*”: i.e., even if $padStr$ was incorrect, it might be assigned a string which accidentally has the correct length, so that the $padLen_5$ variable is nevertheless correct. In this manner, we associate each vertex in the EPG with a conditional probability distribution (CPD), which describes our belief that it was correctly evaluated, conditioned on the correctness status of its immediate predecessors.

Taken together, these CPDs allow us to treat the EPG as a Bayesian network, where each vertex is a Boolean-valued random variable. We also know from executing the tests that some variables in the EPG were incorrect as the test failed, and that some others were (presumably) correct as the test succeeded. This allows us to run the Bayesian network in “*reverse*” and recover the marginal probability $\Pr(\neg l \mid e)$ that each statement l is buggy, conditioned on our observations e on test behaviors. We show the illustrative results of these calculations in Table 1.

Table 1. The ranked list of suspiciousness scores reported by PROSECUTOR for the 10 most suspicious program statements in Figure 1. This is obtained by calculating posterior probability $\Pr(\neg l \mid \neg t_{16})$ for each statement l , given the EPG in Figure 2 as a Bayesian network and the observation $\neg t_{16}$ on the failing execution.

Rank	Suspicion	Line:Statement	Rank	Suspicion	Line:Statement
1	0.531819	26:t1 = padChars[i % padLen]	6	0.166703	12:pads = size - strLen
2	0.261497	25:if (i < pads)	7	0.162269	21:if (pads < padLen)
3	0.205962	9:padChars = str.toCharArray()	8	0.153250	3:if (str == null)
4	0.172209	25:i++	9	0.152615	10:padLen = padStr.length()
5	0.168113	25:i = 0	10	0.150381	11:strLen = str.length()

These suspiciousness scores allow us to naturally rank individual program statements for triage by developers. For the program of Figure 1, one needs to inspect just 3 lines of code to discover the true fault location. More broadly, in our evaluation in Section 5, we find that developers would discover 40% of bugs if they similarly examined only the first 3 statements reported by PROSECUTOR.

2.3 Counterfactual Executions

Consider a second scenario, as shown in Figure 3, where the bug now exists in Line 22. This is a common pitfall: programmers may mistakenly use the wrong method when multiple overloaded methods share the same name. Here, the previously passing test, `testCase1`, now fails since an incorrect string "L" is returned by the statement `padStr.substring(pads)` in Line 22. Notice that during the execution of `testCase1`, if we flip the condition value in Line 21 at Step 11, then the faulty statement would not be executed and, followed by executing Lines 24–28, the test passes successfully. This may provide evidence about either the incorrectness of the condition value (`pads < padLen`)₁₁ in the original failing execution or its dependent variables/values. We can thus leverage such observations by incorporating them into the ranking process.

Similar to the graph in Figure 2, the EPG can be constructed for the failing execution obtained from `testCase1` in Figure 3. Given that an assertion violation occurs on Line 32, by performing marginal inference on the constructed graph, we calculate the posterior probabilities $\Pr(\neg t_i \mid \neg \text{assert}_{13})$ for each branch condition t_i in EPG (blue nodes in Figure 2). This procedure creates a ranked list of branch conditions that are suspected to be incorrectly evaluated. Next, for each condition in the top 20 entries of the ranking, we artificially flip it at runtime and observe its downstream effect on the test outcome. If a previously failing test now passes in this counterfactual scenario, then it might be the case that in the original failing execution (a) this branch condition was evaluated incorrectly, or (b) a value/variable was incorrectly computed under the control of this branch condition. E.g., flipping the branch predicate (`pads < padLen`) at Step 11 changes the test behavior since the value calculated and returned at Step 12 was erroneous.

In such cases, we create a new virtual node, named `cfx`, in the EPG with incoming edges from branch values, flipping of which changes the test outcome and their immediate successors. Finally, we calculate the posterior probability $\Pr(\neg l \mid e \wedge \neg \text{cfx})$ for each statement l as its suspiciousness

```

2 public static String leftPad(String str, int size, String padStr) {
3     ...
9     char[] padChars = padStr.toCharArray(); // fixed
10    ...
21   } else if (pads < padLen) {
22       return padStr.substring(pads).concat(str); // faulty line
23   }
24   char[] padding = new char[pads];
25   for (int i = 0; i < pads; i++) {
26       padding[i] = padChars[i % padLen];
27   }
28   return new String(padding).concat(str);
}

31 public void testCase1() { // Now, an example of failing test
32     assertEquals("!aabc", StringUtils.leftPad("abc", 5, "!aL"));
33 }

```

Fig. 3. A variant of the code in Figure 1 with a bug on Line 22, which should instead call `padStr.substring(0, pads)`.

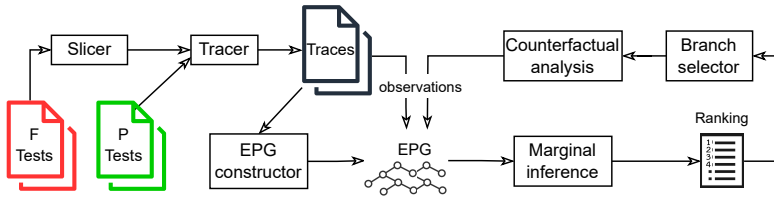


Fig. 4. The PROSECUTOR fault localization workflow.

score. Note that treating cfx as an erroneous variable enables the model to consider that either nodes with incoming edges to cfx could be incorrect, and this suspicion will be further propagated to all other nodes reachable from cfx .

Concretely, feeding the the above counterfactual results back into the model will raises the suspiciousness score of Line 22 from 0.47 to 0.61 and places this line at the top of the ranking.

3 Methodology

Figure 4 presents a high-level overview of the PROSECUTOR fault localization workflow. We now describe this workflow in detail, including the EPG construction process (Section 3.1), a recap of Bayesian networks and the manner in which we probabilistically view the EPG (Section 3.2). Finally, the process of collecting observations, performing counterfactual analysis and computing suspiciousness scores will be elaborated (Section 3.3).

3.1 Constructing the Error Propagation Graph

The EPG is fundamentally constructed from runtime control and data dependencies. We precompute the set of static control dependencies along with DEF/USE sets through a static analysis phase.

3.1.1 Static Program Dependence Analysis. To compute data dependencies in object-oriented languages like Java, we need to consider objects as complex variables with state that can change through manipulation of their attributes. Consider a field-insensitive analysis where changing any attribute of an object changes the entire state of the object. This is straightforward when restricted to the scope of a procedure. One can add the object to the DEF set upon encountering a statement which updates one of its fields. However, it becomes tricky to track object modifications across procedure boundaries and propagate their effects on data dependencies to the caller.

Instead of being fully precise in tracking object modifications, we efficiently approximate the set of defined objects at call sites. We consider that the state of an object u may change when (a) one of the methods of u is called, or (b) a reference to u is passed as an argument to another procedure. We therefore suggest adding receiver objects to the DEF set at call sites as well as the objects and other reference variables, including arrays, which are passed as arguments in call statements. This allows the local def-use chain computation to assume that the new versions of these variables are used in the caller after returning from the callee.

This over-approximation of defined objects might be inefficient due to introducing several spurious data dependencies. Even though invoking a method of an object can change its attributes and therefore the object's state, we assert that this possibility is not identical across different invocations. Invoking *setter-like* methods, which do not return values, including constructors, is very likely to change the state of an object, while *getter-like* methods are often designed to perform computations and eventually return the results. We thus make a more conservative approximation by adding receiver objects to the DEF set only in the case of calling methods with `void` return type. Based on this view of data dependency, we specify DEF/USE sets for primitive statements in Table 2.

Table 2. Defined and used variables in 3-address statements. R_a is the set of reference variables passed as arguments to the callee procedure. The main difference with classical DEF/USE analysis is in method calls, where we treat both receiver objects and reference arguments as being possibly defined.

Stmt type	Stmt format	DEF	USE
Assign Stmt	$t = \text{new Class}(a_1, \dots, a_n)$	$\{t\} \cup R_a$	$\{a_1, \dots, a_n\}$
	$t = \text{new Type}[x]$	$\{t\}$	$\{x\}$
	$t = u.\text{field}$	$\{t\}$	$\{u\}$
	$u.\text{field} = t$	$\{u\}$	$\{u, t\}$
	$t = u[x]$	$\{t\}$	$\{u, x\}$
Return Stmt	$u[x] = t$	$\{u\}$	$\{u, x, t\}$
	$t = u$	$\{t\}$	$\{u\}$
	$t = u_1 \text{ binop } u_2$	$\{t\}$	$\{u_1, u_2\}$
	return t	\emptyset	$\{t\}$
	return	\emptyset	\emptyset
If Stmt	if u_1 binop u_2 goto L	\emptyset	$\{u_1, u_2\}$
Goto Stmt	goto L	\emptyset	\emptyset
Throw Stmt	throw t	\emptyset	$\{t\}$
Catch Stmt	catch e	$\{e\}$	\emptyset
Invoke Stmt	$t = u.\text{method}(a_1, \dots, a_n)$	$\{t\} \cup R_a$	$\{a_1, \dots, a_n, u\}$
	$u.\text{method}(a_1, \dots, a_n)$	$\{u\} \cup R_a$	$\{a_1, \dots, a_n, u\}$
	$t = \text{function}(a_1, \dots, a_n)$	$\{t\} \cup R_a$	$\{a_1, \dots, a_n\}$
	$\text{function}(a_1, \dots, a_n)$	R_a	$\{a_1, \dots, a_n\}$

To obtain control dependencies, we statically calculate dominance frontiers on the reverse control-flow graph (CFG) of each procedure (RDF relations), as proposed by Cytron et al. [11]. This approach determines control dependencies introduced by conventional control-flow constructs (e.g., **if-else**) with explicit edges in CFG. On the other hand, Java programs also make frequent use of **try**, **catch**, and **throw** constructs, which introduce implicit edges from **throw** sites to **catch** statements. For instance, in the code fragment **try** $\{\dots\}$ **catch**(Exception e){ \dots }, several statements inside the **try** block may potentially throw an exception, which would cause control to jump directly to the **catch** statement. We thus assume that every statement in a **catch** block is statically control-dependent on the corresponding **catch** statement, which itself is control-dependent on all statements in the associated **try** block.

3.1.2 Online Graph Construction. We now explain the online process of computing runtime dependencies and constructing the EPG. To specify the EPG, we need to define the set of its vertices and edges. Consider a given trace $\tau = t_1, t_2, \dots, t_n$ in which every statement $t_i \in \tau$ is an instance of a program statement, $l_i = \text{STMT}(t_i)$.

Vertices, V . We include all these program statements in the set of *statement* vertices: $V_l = \{l_i \mid t_i \in \tau\}$. Next, we associate every statement $t_i \in \tau$ with a set of vertices V_i . For every variable $v \in \text{DEF}(l_i)$ defined at t_i , we create a *variable* vertex $\langle v@i \rangle$. If t_i is an object creation or invocation statement, then we also create a variable vertex $\langle x@i \rangle$ for each formal argument, $x \in \text{FARGS}(t_i)$ of its callee function. On the other hand, if t_i does not define any variables, i.e., if $\text{DEF}(l_i) = \emptyset$, then we create a dummy *value* vertex, $\langle \#@i \rangle$. We put these together by defining $V_i = V_i^d \cup V_i^a$, where

$$V_i^d = \begin{cases} \{\langle \#@i \rangle\} & \text{if } \text{DEF}(l_i) = \emptyset, \text{ and} \\ \{\langle v@i \rangle \mid v \in \text{DEF}(l_i)\} & \text{otherwise,} \end{cases}$$

and $V_i^a = \{\langle x@i \rangle \mid x \in \text{FARGS}(t_i)\} \cup \{\langle \text{this}@i \rangle\}$, respectively. Intuitively, V_i^d corresponds to the variable definitions in the procedure where t_i itself is executed, while V_i^a corresponds to definitions of callee arguments created by the call statement t_i upon invocation. We collect the set of value

$$\begin{array}{l}
R_l \frac{v \in V_i}{l_i \rightarrow v \in E} \quad R_{DU} \frac{v \in V_i^d \quad u \in \text{USE}(l_i) \quad \langle u@j \rangle = \text{LASTDEF}(u, t_i)}{\langle u@j \rangle \rightarrow v \in E} \quad R_{CTR1} \frac{v \in V_i \quad t_j = \text{LASTCTRL}(t_i) \quad u \in V_j^d}{u \rightarrow v \in E} \\
R_{CTR2} \frac{v \in V_i^d \quad u \in \text{USE}(l_i) \quad \langle u@j \rangle = \text{LASTDEF}(u, t_i) \quad j < r < i \quad l \in \text{RDFINVERSE}(l_r) \quad u \in \text{DEF}(l)}{\langle \#@r \rangle \rightarrow v \in E}
\end{array}$$

Fig. 5. Inference rules for deriving intra-procedural EPG edges.

and variable vertices: $V_\tau = \bigcup_i V_i$, and define the final set of EPG vertices as:

$$V = V_l \cup V_\tau. \quad (1)$$

Intra-procedural edges. We now define the set of edges E in the EPG. In our implementation, we construct E by making one forward pass over the trace τ . For ease of presentation, we instead use the set of inference rules shown in Figures 5, 7, and 8 to construct E .

The Rule R_l captures the conceptually simplest set of edges: For every trace statement t_i , we include an edge from l_i to each vertex $v \in V_i$. This indicates that the correctness of v depends on the correctness of l_i . Notice that the correctness of vertices $v \in V_i^d$ also depends on the correctness of the operands u that are used. Let c_k refer to the context of a procedure invocation which executes the statement t_k . We write $\text{LASTDEF}(u, t_i)$ to refer to the vertex corresponding to the most recent definition of u that reaches the statement t_i within its context c_i . Rule R_{DU} now creates an edge from $\langle u@j \rangle$ to v , where t_j is the statement creating this variable vertex. These edges correspond to the intra-procedural def-use chains shown in yellow in Figure 2.

Finally, the correctness of each vertex $v \in V_i$ also depends on whether control should have reached this point. At each step i in the trace τ , let $\text{LASTCTRL}(t_i)$ be the most recent statement t_j executed in the same context as t_i and which controls the execution of t_i , i.e., l_i is control-dependent on l_j . In this circumstance, Rule R_{CTR1} enables us to create an edge to v from u for each $u \in V_j^d$. These edges correspond to the intra-procedural control dependence edges shown in blue in Figure 2. Also, note that—regardless of our static over-approximation—this rule only creates edges from the actual location of the `throw` to the corresponding `catch` statement.

Now, recall the assignment `padLen = padStr.length()` in Step 5 (Line 10) of the failing trace in Figure 1. Clearly, this statement is not control dependent on the `if`-statement on Line 6, i.e., `if(padStr.length()== 0)`. Nevertheless, if Line 6 was erroneous, then it might have led to a failure to update `padStr` on Line 7, thereby resulting in an incorrect operand being used in Step 5 of the failing trace. Abstractly, the classical definition of control dependency assumes that a statement l is control-dependent on another statement l' if the outcome of l' determines whether l is executed. Such definition however does not consider the impact of this outcome on the values on l . To account for the effect of such unexplored branches, we introduce new dependency edges into the EPG using Rule R_{CTR2} .

Let $\text{RDFINVERSE}(l_r)$ be the set of all statements that are statically control-dependent on the statement l_r . Consider a statement t_i and one of its operands u . Say that u was last defined at statement t_j in the trace. Now consider all branching statements t_r , executed between t_j and t_i , such that t_r controls the execution of a fourth statement l which also defines u . In this case, Rule R_{CTR2} adds an edge from $\langle \#@r \rangle$ to v .

Example 1. Notice that in Figure 1, the value of the `padStr` operand in Step 5 comes from its definition in the `leftPad` entry point, specified upon invocation at Step 1. Also, the `if`-statement on Line 6 is executed at Step 3, between the definition and the use location of `padStr`. Since the defining statement `padStr = ""` on Line 7 is control-dependent on `if(padStr.length()== 0)` at Line 6, Rule R_{CTR2} derives an edge between nodes $(t0==0)_3$ and `padLen`₄ in Figure 2.

$$\begin{array}{l}
R_{OP1} \frac{OPEN_i \quad u \in USE(l_i) \quad ISARG_i(u) \quad \langle u@j \rangle = LASTDEF(u, t_i) \quad w \in MAP(u, l_i)}{\langle u@j \rangle \rightarrow \langle w@i \rangle \in E} \\
R_{OP2} \frac{OPEN_i \quad u \in USE(l_i) \quad ISRCVOBJ_i(u) \quad \langle u@j \rangle = LASTDEF(u, t_i)}{\langle u@j \rangle \rightarrow \langle this@i \rangle \in E} \\
R_{OP3} \frac{OPEN_i \quad u \in USE(l_i) \quad ISARG_i(u) \quad \langle u@j \rangle = LASTDEF(u, t_i) \quad j < r < i \quad l \in RDFINVERSE(l_r) \quad u \in DEF(l) \quad w \in MAP(u, l_i)}{\langle \#@r \rangle \rightarrow \langle w@i \rangle \in E} \\
R_{OP4} \frac{OPEN_i \quad u \in USE(l_i) \quad ISRCVOBJ_i(u) \quad \langle u@j \rangle = LASTDEF(u, t_i) \quad j < r < i \quad l \in RDFINVERSE(l_r) \quad u \in DEF(l)}{\langle \#@r \rangle \rightarrow \langle this@i \rangle \in E}
\end{array}$$

Fig. 7. Inference rules for deriving inter-procedural EPG edges from call-sites.

The same challenge arises when processing dependencies within syntactically infinite loops. Consider the program in Figure 6 with a `while(true)` loop. Since the `return` statement on Line 6 post-dominates the `if`-statement on Line 5, the former is not statically control-dependent on the latter. However, the correctness of the returned variable `sum` in Step 7 depends on whether the condition in Step 6 was evaluated correctly. To address this challenge, we insert a virtual exit within these loops, such as in Figure 6. Although the `throw` statement, inserted by instrumentation, is unreachable at runtime, it makes a path from Line 5 to a second exit point in the static CFG, and tricks the algorithm into establishing a control dependency between Line 5 and Line 6.

Inter-procedural edges. So far, our model considers correctness dependencies for each procedure in isolation. This allows us to track the propagation of errors across different program elements within each context. Nevertheless, further *inter-procedural* dependencies are required to enable tracking errors across procedure boundaries. We generally need to specify: (a) how the correctness of elements in the caller affects the ones in the callee at the invocation site, and conversely (b) how the correctness of elements in the callee may affect the caller variables/values at the return site. We discuss these rules in Figures 7 and 8 respectively.

Rules R_{OP1} and R_{OP2} involve deriving inter-procedural data dependence edges that specify (a) in EPG. Let $OPEN_i$ be a boolean variable denoting whether the execution of the statement t_i in the trace is followed by the start of a new context. Also, let $ISARG_i(u)$ and $ISRCVOBJ_i(u)$ denote whether the variable u is passed as an argument or the receiver object in statement t_i , respectively. Now, for each operand u , which is last defined by the statement t_j and is passed as the argument w in t_i , Rule R_{OP1} derives an edge from $\langle u@j \rangle$ to $\langle w@i \rangle$. A similar edge is created by Rule R_{OP2} for the receiver object u .

Rules R_{OP3} and R_{OP4} are the inter-procedural counterparts of R_{CTR2} and capture the effect of unexplored branches on the correctness of argument vertices.

Note 3.1. Notice that the previous Rule R_{CTR1} also derives another set of inter-procedural control dependence edges from $u \in V_j^d$ to argument vertices $v \in V_i^a$ created upon invocation, where t_j is the statement which controls the execution of call statement t_i .

The rules in Figure 7 alongside R_{CTR1} therefore fully capture inter-procedural control and data dependencies which are necessary to specify (a).

```

1 int getSum(int [] arr, int size) {
2   int counter = size, sum = 0;
3   while (true) {
4     counter--;
5     if (counter == 0) // faulty
6       return sum;
7     sum += arr[counter];
8     ... instrumentation
9   }
10  if (1 == 0) throw an exception;
11  public void failingTest() {
12    int [] arr = {3};
13    int sum = getSum(arr, 1);
14    assertEquals(3, sum);
15  }

```

Fig. 6. An example program with a bug on Line 5, which should instead be `if (counter < 0)`.

$$\begin{array}{l}
491 \quad R_{CL1} \frac{\text{CLOSE}_i \quad t_j = \text{STKTOP}_i \quad v \in V_i^d \quad u \in V_j^d}{v \rightarrow u \in E} \\
492 \\
493 \quad R_{CL2} \frac{\text{CLOSE}_i \quad t_r = \text{STKTOP}_i \quad v \in \text{DEF}(l_r) \quad \text{ISARG}_r(v) \quad w \in \text{MAP}(v, l_r) \quad \langle w@j \rangle = \text{LASTDEF}(w, t_i) \quad r \neq j}{\langle w@j \rangle \rightarrow \langle v@r \rangle \in E} \\
494 \\
495 \quad R_{CL3} \frac{\text{CLOSE}_i \quad t_r = \text{STKTOP}_i \quad v \in \text{DEF}(l_r) \quad \text{ISRCVOBJ}_r(v) \quad \langle \text{this}@j \rangle = \text{LASTDEF}(\text{this}, t_i) \quad r \neq j}{\langle \text{this}@j \rangle \rightarrow \langle v@r \rangle \in E} \\
496 \\
497 \\
498 \\
499
\end{array}$$

Fig. 8. Inference rules for deriving inter-procedural EPG edges from return locations.

On the other hand, inter-procedural edges that specify (b) are derived by Rules R_{CL1} – R_{CL3} . These rules activate when a context is closed. To calculate these edges, we need to first process all the executed statements within the callee context. Let the boolean variable CLOSE_i denote whether the statement t_i is the last executed statement within the context c_i . Usually, the context closes correspond to return statements. Note, however, that they capture a broader class of program statements, including crash-inducing and exception-throwing ones. Note also that when a function invocation does not gracefully return, the call statement itself serves as the point of context close. Also, let STKTOP_i denote the latest statement t_j executed in the trace before starting the context c_i . These variables allow us to link two consecutive call frames/contexts within the call stack.

Recall once again the discussion of how we encode the potential impact of unexplored branches as dependency edges within the EPG. We apply a similar reasoning to the point where the context is closed: In particular, context closing events (e.g., return) might prevent subsequent desirable updates from being applied to reference arguments of the caller, or they might return incorrect values which affect the correctness of variables storing them in the caller. Rule R_{CL1} thus preemptively creates dependency edges $v \rightarrow u$ from the return statement v to all variable vertices u created at the call site. Let t_r be the call statement. Rule R_{CL2} creates an edge from $\langle w@j \rangle$ to $\langle v@r \rangle$ for each reference variable v that is passed as an argument w , and where t_j is the statement which last defined w in the callee context. Rule R_{CL3} defines similar edges for the receiver object.

3.2 A Probabilistic View of the EPG

Conditional probability distributions. Each vertex in the EPG, $G = (V, E)$, represents the correctness of some program element, i.e., either a statement, control flow decision, or a runtime value. As previously discussed in Section 2, the edge set E encodes possible ways in which errors may propagate along the program execution.

The following principle is fundamental to the PROSECUTOR workflow: If all of the predecessors u of a node v in the EPG are correctly evaluated, then v is itself correctly evaluated. Conversely, if any of its predecessors u is incorrect, then v *might* itself be incorrect.

We quantify this idea by treating each program element $v \in V$ of the EPG as a Boolean-valued random variable and associating it with a conditional probability distribution (CPD). We write $\text{Pred}(v) = \{u \in V \mid u \rightarrow v \in E\}$ for the set of predecessors of v in the EPG. The conditional probability distribution of v is a function P which maps a valuation $\mathbf{x}_{\text{Pred}(v)}$ of the predecessors of v to the probability that v is true (“*correctly evaluated*”). We define:

$$P(v \mid \mathbf{x}_{\text{Pred}(v)}) = p^w, \quad (2)$$

where $p = 0.85$ and $w = \#\text{false}(\mathbf{x}_{\text{Pred}(v)})$ is the number of predecessors u that were incorrectly evaluated. Conditional probabilities must add up to 1, so we define $P(\neg v \mid \mathbf{x}_{\text{Pred}(v)}) = 1 - p^w$.

Notice that when all predecessors were correctly evaluated, $w = 0$, so that $P(v \mid \mathbf{x}_{\text{Pred}(v)}) = 1$. Notice also that the conditional probability of v drops as w increases (but never reaches 0, to account

for the possibility of fault-tolerant program statements). For vertices v with no incoming edges in the EPG (e.g., statement vertices $l_i \in V_l$), we uniformly use the prior probability, $\Pr(v) = 0.85$.

SmartFL associates each type of statement (arithmetic operations, function calls, etc.) with a manually chosen probability of propagating errors. In contrast, observe that PROSECUTOR considers all statements to be equally fault-tolerant, and instead sets up density functions so that the correctness probability of the output is dependent on the number of erroneous inputs.

Bayesian networks. Note that our ultimate goal is to calculate posterior suspicions, $\Pr(\neg l \mid e)$, that each program statement l is incorrect, given our observations e on the test outcomes. In order to speak meaningfully about these conditional probabilities, we need to set up a joint probability distribution over the variables of V .

At this point, we note that, by construction, the EPG is an acyclic graph. As a result, together with the CPDs $P(v \mid \mathbf{x}_{\text{Pred}(v)})$ just defined, it forms a Bayesian network [19]. Accordingly, let \mathbf{x} be a valuation of each variable in V . We define:

$$\Pr(\mathbf{x}) = \prod_v P(x_v \mid \mathbf{x}_{\text{Pred}(v)}). \quad (3)$$

It is clearly the case that $\Pr(\mathbf{x}) \geq 0$ and it can be verified that $\sum_{\mathbf{x}} \Pr(\mathbf{x}) = 1$. It follows that $\Pr(\mathbf{x})$ is a well-defined probability distribution. With this definition in hand, one can readily speak of the probability of individual events, $\Pr(e) = \sum_{\mathbf{x} \sim e} \Pr(\mathbf{x})$, as the sum of the probabilities of matching valuations e , and of conditional probabilities, $\Pr(e_1 \mid e_2) = \Pr(e_1 \wedge e_2) / \Pr(e_2)$. In practice, these can be computed by off-the-shelf engines such as LibDAI [23].

3.3 Evidence Collection

To calculate posterior suspicion of individual elements, we collect evidence about the correctness/incorrectness of runtime values/variables from: (a) Passing and failing executions on the original program, e.g., a failed assertion shows the expression inside was evaluated to false, and (b) passing executions on *mutated* versions of the program, e.g., if flipping a branch predicate causes a failing test to succeed, then this condition was likely to be evaluated incorrectly in the original execution.

We consider the trace $\tau = t_1, t_2, \dots, t_n$ obtained upon running each test case. In the case of a failing execution, we mark variables and values $v \in V_n^d$ created at the *failure point* of the trace t_n as being incorrect. As discussed earlier, when a function does not return gracefully (mostly due to a program crash), the call statement t_i in its parent context is treated as the point of context close. In this situation, the variable node which stores the return value at t_i must be intuitively considered incorrect. We thus additionally mark all variables $\{v \in V_i^d \mid \text{CLOSE}_i \wedge \text{OPEN}_i\}$ as incorrect.

In passing executions, on the other hand, we mark *all* variables and values created throughout the trace, $v \in \bigcup_i V_i$ as being correctly evaluated. Although this would lower the suspiciousness $\Pr(\neg l \mid e)$ of program statements l executed in passing traces, recall that each statement still has a non-zero chance of being buggy, because of the fault tolerance built into our choice of CPDs. Furthermore, we assume that all the statements $\{l_i = \text{STMT}(t_i) \mid t_i \in \tau\}$ inside the test procedure are bug-free and mark these statement vertices as being correct.

Note 3.2. A more precise strategy may restrict correctness labels only to variables that appear in passing assertions. However, many real-world tests are *assertion-free*. In fact, across our entire benchmark dataset, 36% of the test cases do not contain assertions. In such cases, the absence of a failure during test execution is implicitly interpreted as a “pass”. To account for such tests, we uniformly consider all variables and values created during passing executions as being correct. Another more principled approach might involve asking the user to inspect and label the correctness of intermediate values by hand. This approach would be similar to that employed by Zhang et al. in Ursa [39] and Xu et al. in [36]. The main challenges that we anticipate in this approach would be

589 in determining which values to present to the user for inspection, in providing them with enough
 590 context to make their judgments, and in gracefully permitting the possibility of user error. We
 591 believe that this is an exciting direction for future research.

592 *Counterfactual analysis.* We augment our evidence about the incorrectness of program elements
 593 by performing a set of counterfactual experiments. The key idea is that if changing a value at
 594 runtime can make a failing test succeed, then this value is likely to be related to the bug. In particular,
 595 we focus on changing branch conditions as their domain has only two values, `false` and `true`.

596 After running the original failing test cases, we perform an initial Bayesian inference run to
 597 identify the $k = 20$ most suspicious runtime branch conditions. For each of these branch conditions
 598 t_r , in sequence: We surgically alter the program to flip the value of each specific branch condition
 599 at runtime, and rerun the failing test to determine whether this alteration causes the test to instead
 600 succeed. If the test outcome is changed, it might be the case that in the original trace: (a) this branch
 601 condition—specifically the value vertex $\langle \# @ r \rangle$ —was evaluated incorrectly, or (b) this condition
 602 results in executing a statement t_k that introduces or propagates an error into the program state.
 603 We account for both possibilities by introducing a new vertex `cfx` into the EPG. We also add the
 604 following edges to E : (a) $\langle \# @ r \rangle \rightarrow \text{cfx}$, and (b) $v \rightarrow \text{cfx}$ for each vertex v that is an immediate
 605 successor of $\langle \# @ r \rangle$. We then mark the virtual variable `cfx` as being incorrectly evaluated.

606 There are two principal advantages of our approach over traditional mutation-based approaches:
 607 First, MBFL does not consider the effect of modifying a statement on its dependent statements. In
 608 other words, they only alter suspiciousness of the mutated statement in response to a successful
 609 mutation. Second, MBFL is not equipped with a built-in heuristic to restrict the large space of
 610 mutations. This imposes substantial overhead to run hundreds of experiments, which requires a
 611 huge time budget and drastically restricts the scalability of these techniques.

613 4 Implementation

614 We have implemented PROSECUTOR using the Soot framework [30], which translates Java bytecode
 615 into typed 3-address Jimple code. Our implementation consists of ≈ 3000 lines of Java to perform
 616 instrumentation and static analysis, and 2200 lines of Python code to construct the EPG. We also
 617 use LibDAI for marginal inference [23].

618 As shown in Figure 4, PROSECUTOR captures the last statement t_n executed in each failing test
 619 and extracts an executable backward slice, $\bigcup_{u \in \text{USE}(l_n)} \text{BACKSLICE}\langle l_n, u \rangle$, which transitively records
 620 all statements affecting a used variable u at statement l_n [9]. This allows us to have near-minimal
 621 failure-inducing tests, which reduce the size of traces by preemptively eliminating several useless
 622 paths in EPG. To construct the EPG, we use all traces F obtained from the sliced versions of the
 623 failing tests. In order to speed up EPG construction and reduce the size of the final graph, we
 624 perform a coarse-grained coverage analysis, and heuristically choose a set of passing traces P that
 625 have large overlap in their executed methods with F .

626 *Rank aggregation.* Intuitively, the traces in P provide statistical coverage information similar
 627 to SBFL. Although this information is useful, sometimes passing traces frequently execute faulty
 628 statements without themselves failing. This leads to cases where the suspiciousness of faulty
 629 statements is incorrectly suppressed, in a manner similar to the shortcomings of SBFL. We therefore
 630 construct the EPG and perform marginal inference in two settings: (a) using only failing traces in F ,
 631 and (b) using traces in $F \cup P$. As we will see in Section 5, running PROSECUTOR in mode F already
 632 outperforms all baselines, and $F \cup P$ leads to additional improvements. Ultimately, PROSECUTOR
 633 combines the two rankings, F and $F \cup P$, by ranking statements according to $r_l = \min(r_l^{F \cup P}, r_l^F)$ and
 634 by giving priority to $F \cup P$ to break ties. This enables us to integrate two ranking perspectives: F
 635 solely considers the error propagation across failing traces, and $F \cup P$ partially considers coverage
 636

Algorithm 1: COMPRESS(τ, γ), where $\tau = \langle t_1, t_2, \dots, t_n \rangle$ is a trace, and $\gamma > 1$ is the minimum desired repetition period.

1. For each program statement l that is a loop entry point, define:

$$\chi(l) = [i \mid \text{STMT}(t_i) = l].$$

2. Over all loop entry points l , in decreasing order of number of occurrences, $|\chi(l)|$:

A. Let $\chi(l) = [i_1, i_2, \dots, i_k]$.

B. For each $i_a \in \chi(l)$, if ISLOOP($\tau_{i_a}, \dots, \tau_{i_a+\gamma}$):

a. Find the largest $b \geq a + \gamma$ such that ISLOOP($\tau_{i_a}, \dots, \tau_{i_b}$). ▷ Can be done through Binary Search

b. Let L be the loop period reported by ISLOOP.

c. Let $\sigma = t_{i_b-2L}, \dots, t_{i_b-1}$ and $\tau' = t_1, t_2, \dots, t_{i_a-1}, \sigma, t_{i_b}, \dots, t_n$. ▷ Keep 2 iterations of the loop

d. Return COMPRESS(τ', γ).

3. If $\gamma > 3$, return COMPRESS($\tau, \gamma - 1$). ▷ No more loops with minimum period γ could be found

4. Return τ . No more loops could be found.

of passing traces besides the propagation of errors. We will show in our experimental evaluation that this ranking aggregation achieves better performance than either approach alone.

Trace compression. Many test cases in the Defects4J suite execute long running loops, which increases the overhead of EPG construction and slows down marginal inference. On the other hand, these long running loops produce regular and repetitive data and control dependencies. This allows us to eliminate repeated loop iterations to keep the EPG compact, and accelerate both its construction and subsequent Bayesian inference.

Given a trace $\tau = t_1, t_2, \dots, t_n$, a loop is identified by a subsequence of statements $t_i, t_{i+1}, \dots, t_{i+\ell}$ repeated consecutively in τ . To preserve data and control dependencies within a loop, it suffices to retain two iterations to capture dependencies both upon *entering* and *exiting* the loop. In this case, data and control dependencies from the remaining iterations are identical to the last iteration. We use the instrumentation-guided loop identification algorithm shown in Algorithm 1 to find and compress loop iterations. Note that ISLOOP is a decision algorithm with time complexity $O(n)$, straightforward to implement using dynamic programming and KMP.

When instrumenting programs at the IR level in Soot, we identify `goto` statements l_g along with their corresponding target statements l_t , and selectively insert annotations immediately before l_t if $l_t.\text{lineNumber} \leq l_g.\text{lineNumber}$. This represents a jump within a procedure to one of its preceding statements, which naturally represents a backedge in the CFG. We use this instrumentation to identify statements that serve as loop entry points in a given trace.

Aliasing. Another challenge arises when multiple variables point to the same memory location and defining a variable u might result in an “*action-at-a-distance*” update to another variable v . A thorough solution to this problem would require careful modeling of the heap. We instead perform a lightweight analysis to create the set REF of all assignment statements whose right-hand side is a reference variable. Next, when computing DEF(l), we consult REF to determine if there is an aliasing variable w for each $v \in \text{DEF}(l)$. In this case, we add w to DEF(l) and transitively perform this augmentation until fixpoint. Note that we assume an aliasing remains persistent through the execution, which is guaranteed by the SSA-like format that Soot provides.

Pruning dependencies. Finally, recall that the conditional probability distributions defined in Equation 2 contain 2^k entries, where $k = |\text{Pred}(v)|$ is the in-degree of the vertex v . In particular, this exponential growth of the CPD tables causes challenges with: (a) invocation statements with several arguments, and with (b) virtual variables `cfx` that may depend on many vertices.

To avoid space explosion in (a), in the case of many-argument functions (i.e., where $|\text{FARGS}(t_i)| > 10$), we delete edges from argument definitions to $v \in V_i$. This optimization may be alternatively

characterized as disabling Rule R_{DU} for many-argument invocation statements. Note that these dependencies can still be recovered by transitively following inter-procedural edges, so the optimization does not lead to a total loss of error propagation pathways.

To address challenge (b), arising from cfx variables with high in-degree, we break dependencies with introducing new variables cfx_1, \dots, cfx_n , each of which is dependent on at most 10 other vertices and which are themselves connected to the main virtual variable cfx . This limits the size of the CPDs to 2^{10} entries, while leading to at most a linear increase in the number of EPG vertices.

5 Experimental Evaluation

Our evaluation seeks to answer the following questions:

RQ1. How effective is PROSECUTOR in localizing faults?

RQ2. How long does PROSECUTOR take to compute its ranking?

RQ3. How do different aspects of EPG construction affect ranking quality?

RQ4. How do counterfactual executions affect the accuracy?

Experimental setup. We conducted all our experiments on a workstation machine with an AMD Ryzen 9 5950X CPU and 128 GB of memory running Ubuntu 22.04. We set a timeout of 20 hours per faulty version for PROSECUTOR and the baselines.

Benchmarks. We evaluate our approach on 470 faulty versions of 13 projects from the Defects4J (v3.0.1) benchmark suite [14]. These faulty versions contain, on average, 21K lines of source code. In addition, each failing test on average consists of a trace of length 22,000 steps and executes ≈ 400 distinct lines of code overall. Even though loop compression and dependency pruning massively improve the overall scalability of our approach, the current implementation still faces challenges when analyzing recursive programs or those with hundreds of sparse loops (as we need to keep 2 iterations of each loop). In such cases, both EPG construction and the subsequent marginal inference process incur substantial computational overhead. Among the total 542 versions of these 13 projects, we thus exclude benchmarks whose failing traces are longer than 500,000 Jimple statements—there are 61 such buggy versions in all. We also exclude 11 benchmarks with incorrect function declarations or missing method bodies, where no *executable* line can be blamed.

Baselines. We compare PROSECUTOR to 8 existing techniques including:

- (1) four SBFL approaches: Op2 [25], Tarantula [13], DStar [32], and Ochiai [6],
- (2) two MBFL approaches: MUSE [24] and Metallaxis [26], and
- (3) two other state-of-the-art methods: SmartFL [35] and DepGraph [27], respectively.

We reimplemented SBFL and MBFL approaches using Soot analysis framework [30] to get coverage profiles through instrumentation and by using Major [5] to generate tentative mutants. We also used the implementations of SmartFL and DepGraph provided in their available artifacts.

5.1 RQ1: Effectiveness of Ranking

We started by identifying the actual faulty lines in all benchmarks and consulted the patches provided by Defects4J to derive the expected ground truth.

5.1.1 Quantitative Analysis. Next, we ran both our system and the baselines on all benchmarks and noted the position of the *faulty line* in each proposed ranking. We aggregate these results in Table 3, using (a) top- k measure that shows the number of project versions whose bug was successfully identified within the top k entries of the reported rankings, and (b) the median EXAM score, which measures the median percentage of code examined to locate the faulty statement in each ranking.

Table 3. Effectiveness of PROSECUTOR and the baselines in statement-level fault localization. The best performing approach in each case is marked in bold.

Project	Metric	DStar	Ochiai	Op2	Tarantula	Muse	Metallaxis	SmartFL	PROSECUTOR
Lang (58)	TOP-1	7	7	7	8	8	13	15	20
	TOP-3	20	20	20	21	16	26	30	31
	TOP-5	25	26	23	26	22	34	34	41
	TOP-10	39	38	39	38	29	37	42	49
	TOP-20	45	42	45	43	35	45	47	56
	EXAM	14.0	14.1	14.1	14.5	21.2	9.4	9.6	6.0
Compress (45)	TOP-1	10	9	10	8	2	9	6	13
	TOP-3	14	13	14	12	9	17	13	19
	TOP-5	19	17	19	16	12	18	17	21
	TOP-10	22	21	22	21	20	22	23	28
	TOP-20	23	22	23	23	22	26	27	34
	EXAM	14.4	14.8	14.4	14.4	20.0	9.2	10.3	5.0
Chart (26)	TOP-1	3	3	2	2	2	2	5	5
	TOP-3	7	7	6	7	3	5	14	17
	TOP-5	10	10	9	10	4	6	16	21
	TOP-10	11	12	10	12	7	8	18	23
	TOP-20	14	16	15	16	8	10	20	25
	EXAM	8.3	7.6	9.0	8.3	35.4	33.8	5.1	2.0
JackCore (25)	TOP-1	4	4	5	4	1	4	7	6
	TOP-3	8	7	8	8	3	6	10	13
	TOP-5	10	10	10	10	4	11	12	14
	TOP-10	12	11	12	12	6	12	13	17
	TOP-20	13	14	13	15	8	12	14	18
	EXAM	4.0	5.2	3.9	3.5	26.6	3.9	5.9	2.0
Math (82)	TOP-1	15	15	14	14	12	9	14	13
	TOP-3	25	25	23	25	17	20	27	28
	TOP-5	29	29	27	30	22	27	35	38
	TOP-10	35	37	33	37	31	35	39	50
	TOP-20	43	45	41	45	41	40	47	57
	EXAM	9.1	9.1	10.3	8.3	16.9	10.3	8.3	5.0
Jsoup (60)	TOP-1	6	6	5	6	3	4	5	14
	TOP-3	12	12	11	13	4	11	15	20
	TOP-5	18	17	18	16	6	14	18	24
	TOP-10	27	27	25	22	7	20	21	28
	TOP-20	33	34	31	31	10	22	25	34
	EXAM	2.9	2.9	3.2	3.4	100	100	20.7	2.9
Cli (39)	TOP-1	9	8	9	8	6	4	7	7
	TOP-3	13	12	12	11	7	10	11	10
	TOP-5	18	17	17	17	8	13	12	11
	TOP-10	21	21	20	20	9	17	13	18
	TOP-20	24	24	23	24	12	19	16	23
	EXAM	8.3	8.3	10.7	8.3	32.3	12.2	20.6	12.5
Codec (16)	TOP-1	3	3	3	2	3	3	5	5
	TOP-3	5	5	5	5	4	4	7	6
	TOP-5	5	5	5	5	4	5	8	6
	TOP-10	6	6	6	6	6	6	8	10
	TOP-20	11	11	13	11	8	9	10	12
	EXAM	15.1	15.1	13.9	15.1	20.3	14.4	8.2	9.2
Mockito (38)	TOP-1	6	7	6	6	1	2	9	7
	TOP-3	12	13	11	10	3	5	13	14
	TOP-5	16	18	15	16	5	5	14	14
	TOP-10	19	20	17	19	5	7	15	18
	TOP-20	22	22	19	22	6	7	17	25
	EXAM	2.3	2.0	2.6	2.3	100	100	16.5	1.4
Time (26)	TOP-1	2	2	2	2	1	2	4	3
	TOP-3	5	4	5	5	2	5	7	5
	TOP-5	10	9	9	9	3	7	7	8
	TOP-10	12	12	11	11	4	7	10	12
	TOP-20	15	15	15	16	5	10	11	15
	EXAM	1.5	1.5	1.5	0.9	30.0	6.3	13.9	1.3
Csv (16)	TOP-1	4	4	4	3	2	3	4	6
	TOP-3	6	6	6	5	3	5	5	7
	TOP-5	7	7	7	7	4	7	6	7
	TOP-10	8	8	8	8	6	7	6	8
	TOP-20	12	12	12	12	7	8	7	9
	EXAM	8.5	8.5	8.5	8.5	32.5	24.0	23.0	12.2
Gson (17)	TOP-1	4	5	4	5	1	2	8	3
	TOP-3	7	7	7	8	3	3	11	7
	TOP-5	7	7	7	8	3	3	11	11
	TOP-10	13	11	10	11	3	3	12	11
	TOP-20	14	14	11	13	3	3	12	13
	EXAM	2.5	2.5	5.0	3.4	100	100	1.9	1.9
Collections (22)	TOP-1	5	5	5	5	n/a	n/a	6	3
	TOP-3	10	11	11	11	n/a	n/a	9	13
	TOP-5	14	15	14	15	n/a	n/a	11	16
	TOP-10	17	19	17	19	n/a	n/a	14	18
	TOP-20	18	19	17	19	n/a	n/a	18	21
	EXAM	7.0	5.8	6.2	5.8	n/a	n/a	11.1	5.6
Total (#70)	TOP-1	78	78	76	73	42	57	95	106
	TOP-3	144	142	139	141	74	117	172	189
	TOP-5	188	187	180	185	97	150	201	231
	TOP-10	242	243	230	236	133	183	236	289
	TOP-20	287	290	278	290	165	211	271	342
	EXAM	7.1	6.9	7.7	6.8	29.8	16.7	10.2	4.2

Overall, we observe that PROSECUTOR allows for 40% of true fault locations to be identified by examining just 3 lines of code. Alternatively, the median EXAM score reveals that half of all faulty locations can be identified by examining just 4.2% of the statements covered by failing tests. We include an alternative visual presentation of these results for the entire dataset in Figure 9. Observe that the technique significantly outperforms *all* of the baselines, and identifies 15% more bugs within its top 5 predictions than the best-performing one. Figure 9 also demonstrates that to identify fault locations in 50% of the benchmarks, a developer using PROSECUTOR would need to inspect 44% fewer lines than even the best-performing baseline.

Figure 10 analyzes the overlap between our approach and each baseline. Across all comparisons, observe

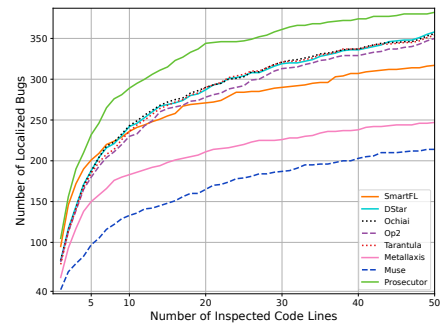


Fig. 9. Number of faults in top- k predictions of each technique across the entire dataset. The ideal technique would place the true fault location as its first prediction, so faster rising curves indicate higher effectiveness.

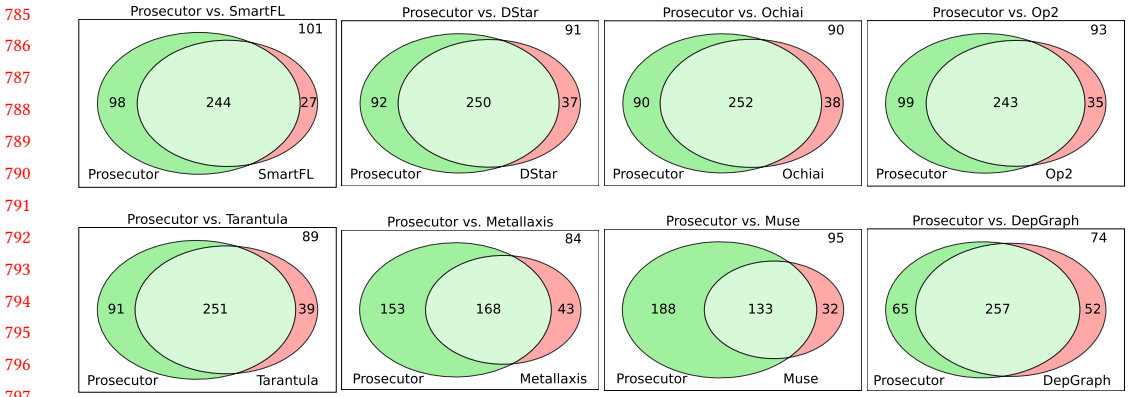


Fig. 10. Venn diagrams illustrating the breakdown of faults localized or missed by each pair of approaches. Each diagram highlights the overlap of faults correctly localized by both methods and incorrectly missed by the two as well as those uniquely localized by each technique. Bugs ranked within the top-20 are considered localized in all diagrams for statement-level localization, while the last chart uses a top-3 threshold for method-level localization. Benchmarks from the Collections project were excluded from the comparison with Muse, Metallaxis and DepGraph.

that PROSECUTOR consistently localizes substantially more bugs than the baselines. Overall, the number of bugs uniquely localized by our approach (the dark green regions) is at least $2.3\times$ greater than the number localized by each individual baseline technique (the red portions of the figure) and $3.6\times$ greater than that of SmartFL. This again highlights the advantage of PROSECUTOR over the baseline techniques.

Our technique demonstrates dramatic improvements on bugs associated with the projects such as Compress and Jsoup. While it is hard to pinpoint an exact cause, we speculate that this is because of challenges faced by SBFL and MBFL techniques when the faulty line is frequently executed by both passing and failing tests. Notice also that SBFL approaches occasionally outperform PROSECUTOR (notably in the Cli project). We believe that this is due to a combination of factors, including unexpectedly deep faults and imperfect modeling in the EPG (such as the lack of aliasing information). Addressing these limitations is a good direction of future work.

Note 5.1. Unfortunately, we could not run MBFL on the Collections project as it needs Java versions not supported by Major.

Comparison to learning-based techniques. We separately present the comparison of PROSECUTOR to DepGraph, a recent state-of-the-art learning based fault localization technique, in Table 4. We separate these results because DepGraph only performs *method-level* localization. In order to perform an apples-to-apples comparison with PROSECUTOR, we coarsen our results by choosing the most suspicious statement within each method as a representative for the method as a whole.

Apart from two projects, Math and Mockito, PROSECUTOR consistently matches or surpasses DepGraph in its top 2, 3, 4, and 5 predictions for the most suspicious methods. Our approach even outperforms DepGraph at the first prediction in several projects, including Lang and Compress, albeit in general DepGraph demonstrates superior prediction for the top 1.

The superior top-1 performance of DepGraph can be reasonably explained by its design focus on method-level fault localization, in contrast to the statement-level suspiciousness scores that are reported by PROSECUTOR. Furthermore, their technique relies on additional sources of information such as the history of code changes, which our approach does not take into account.

Table 4. Comparison of PROSECUTOR and DepGraph in method-level fault localization.

Project	Method	TOP-1	TOP-2	TOP-3	TOP-4	TOP-5
Lang (58)	DepGraph	42	47	51	54	55
	PROSECUTOR	44	55	56	56	56
Compress (45)	DepGraph	23	30	32	35	37
	PROSECUTOR	29	36	39	41	41
Chart (26)	DepGraph	15	16	20	21	21
	PROSECUTOR	12	21	24	24	25
JackCore (25)	DepGraph	10	13	14	14	14
	PROSECUTOR	8	14	17	18	20
Math (82)	DepGraph	51	64	69	75	75
	PROSECUTOR	43	55	61	67	72
Jsoup (60)	DepGraph	17	25	28	29	31
	PROSECUTOR	20	25	32	32	33
Mockito (38)	DepGraph	17	23	26	29	30
	PROSECUTOR	17	22	25	28	28

Project	Method	TOP-1	TOP-2	TOP-3	TOP-4	TOP-5
Codec (16)	DepGraph	7	9	10	10	11
	PROSECUTOR	8	10	11	14	14
Time (26)	DepGraph	14	15	16	16	17
	PROSECUTOR	10	16	16	16	18
Cli (39)	DepGraph	16	21	22	27	28
	PROSECUTOR	10	15	17	21	23
Csv (16)	DepGraph	5	7	7	9	10
	PROSECUTOR	8	9	10	10	11
Gson (17)	DepGraph	13	13	14	14	15
	PROSECUTOR	5	13	14	14	14
Collections (22)	DepGraph	n/a	n/a	n/a	n/a	n/a
	PROSECUTOR	6	14	17	18	20
Total (470)	DepGraph	230	283	309	333	344
	PROSECUTOR	220	305	339	359	376

Note 5.2. (a) We were unable to run DepGraph on the Collections project because of compatibility issues with the most recent version of Defects4J. (b) In contrast to our implementation, DepGraph does not consider the possibility of the bug being present in the initialization portion of the test suite. In order to keep the comparison uniform, our statistics for PROSECUTOR in Table 4 similarly exclude these methods from consideration.

5.1.2 Qualitative Analysis. We also conducted a qualitative analysis to identify scenarios in which PROSECUTOR outperforms SmartFL and vice versa. We now elaborate on such cases through a series of representative examples.

Note that SmartFL creates its probabilistic graph by parsing Java bytecode traces and simulating their execution using an operand stack, where pushing a variable onto the stack is treated as a definition and popping it is treated as a use. In contrast, PROSECUTOR does not rely on stack-based execution to construct the EPG. Instead, our EPG is built upon 3-address statements which are executed along the trace given the inference rules in Figure 5, 7, and 8 and the DEF/USE sets in Table 2. The following examples clarify the fundamental differences in our graph construction.

Example 2. Consider the test case in Figure 11 which fails on Line 7 since the variable p2 is incorrectly computed as a result of a bug within the method combine. When processing the trace and the method call at Step 5, PROSECUTOR derives an edge from the vertex $\langle p2@4 \rangle$ to the vertex $\langle subs@5 \rangle$. Other similar edges in the EPG will be as follows: $\langle p2@1 \rangle \rightarrow \langle p2@4 \rangle$, $\langle p3@2 \rangle \rightarrow \langle p3@3 \rangle$, and $\langle p3@3 \rangle \rightarrow \langle p2@4 \rangle$. These “shortcut” edges are derived by Rule R_{DU} in Figure 5 given the DEF/USE sets defined for invocation statements in Table 2.

In contrast, upon processing the method call on Line 6, SmartFL pushes the value of p2 onto the operand stack and creates a corresponding node. Whenever this value is popped from the stack within the subset method, SmartFL introduces an edge from this vertex to the point of use. The value popped from the stack may be used to define a second variable which itself can influence a third variable and so on. Observe that this chain of dependencies may grow arbitrarily long and involve multiple call frames, eventually propagating to and affecting the return value of the method subset and therefore the value of the subs object.

```

1 public void testCase() {
2   ExtendedProperties p2 = new ExtendedProperties();
3   ExtendedProperties p3 = new ExtendedProperties();
4   p3.setProperty("sub.test", "foo");
5   p2.combine(p3);
6   ExtendedProperties subs = p2.subset("sub");
7   assertNotNull(subs);
8 }
9 public ExtendedProperties subset(String prefix) {
10  ExtendedProperties c = new ExtendedProperties();
11  ... /* several statements inside the method */
12  return c;
13 }

```

Fig. 11. Example test snippet adapted from the Apache Commons Collections library [1] to describe the bug in version #12 of this project in Defects4J. The program contains a bug in the combine method, which results in an incorrect update to the variable p2 on Line 5, and an assertion violation.

883 In practice, an incorrect value which is returned/created as a result of a bug within a method (in
 884 our case, the `combine` method on Line 5) may be passed as an input argument or used as the base
 885 object in subsequent calls, where it can propagate and taint other values. Still, as long as the value
 886 remains structurally valid, subsequent methods may process it without issue (e.g., the method calls
 887 on Lines 6), and the error typically surfaces only when the value becomes invalid or is explicitly
 888 checked by an assertion. The graph representation which is used by SmartFL performs poorly in
 889 such scenarios, where the actual fault resides in higher-level methods far from the failure location.

890 As a consequence, while PROSECUTOR surfaces this bug at Rank 5, its position in SmartFL is lower,
 891 at Rank 10. We perform a more comprehensive quantitative analysis of the effect of these shortcut
 892 edges in our ablation study in Section 5.3.3. In particular, removing shortcut edges in this example
 893 (using the alternative Rule R'_{DU}) causes the actual fault to drop to Rank 9 in the PROSECUTOR ranking.

894 The EPG also introduces an additional set of shortcut edges that more effectively model error
 895 propagation when the failure occurs within a nested callee and results in an abrupt program
 896 termination while multiple frames remained unpopped on the call stack. As discussed earlier,
 897 regardless of whether a method terminates gracefully, Rule R_{CLI} derives edges from vertices defined
 898 by context-closing statements to those in V_j^d , where t_j is the call statement in the parent context.

899 *Example 3.* Consider the program in Figure 12 with
 900 a bug on Line 3 whose execution leads to an unex-
 901 pected invocation of the method `readLog` at Step 3.
 902 When constructing the EPG for this example, our infer-
 903 ence rules specifically derives the following two edges:
 904 (a) Rule R_{CLI} derives an edge from the vertex $\langle \# @ 7 \rangle$
 905 corresponding to the thrown exception to the variable
 906 vertex $\langle \text{logText} @ 3 \rangle$, and (b) Rule R_{CTR1} derives an edge
 907 from the value vertex $\langle \# @ 2 \rangle$ to the vertex $\langle \text{logText} @ 3 \rangle$.

908 Although the variable `logText` is never instantiated
 909 at runtime, we still include the corresponding vertex in
 910 the EPG because our EPG construction operates on 3-
 911 address statements rather than on the concrete stack ex-
 912 ecution. In addition to the value vertex $\langle \# @ 7 \rangle$ at failure
 913 point, we also mark the variable vertex $\langle \text{logText} @ 3 \rangle$
 914 as being incorrect since at a minimum, the method `readLog` must return gracefully for the value
 915 `logText` to have any chance of being correct. The observation $\neg \langle \text{logText} @ 3 \rangle$ along with the edges (a)
 916 and (b) allows the model to assign a high suspiciousness score to the `if`-statement on Line 3.

917 In contrast to PROSECUTOR, SmartFL does not place Line 3 among its top predictions within the
 918 ranking. This occurs since the predicate corresponding to this line is not reachable from the failure
 919 point in the underlying graph and therefore the error does not propagate appropriately. SmartFL's
 920 graph construction assumes that the callee procedure consistently uses the arguments passed at
 921 the call site, and these data dependency edges enable error propagation to top-level procedures. Of
 922 course, if such arguments existed, the corresponding vertices would be pushed onto the stack under
 923 the control of the predicate on Line 3. If these values were later popped within the method `readLog`
 924 and contributed to error propagation, then the statement on Line 3 would be reachable from the
 925 failure location in SmartFL's graph. However, in the case of program in Figure 12, where the callee
 926 method `readLog` receives no arguments, this graph construction fails to connect two consecutively
 927 executed methods, even though the root cause of the fault resides in the caller procedure.

928 The above example simply demonstrates how PROSECUTOR more effectively localizes faults in
 929 higher-level procedures by incorporating shortcut edges that connect consecutive incomplete call
 930
 931

```

1 public static void processLog() {
2   boolean ready = writeLog();
3   if (!ready) // faulty line
4     String logText = readLog();
5   ... /* remaining code comes here */
6 }
7 public static String readLog() {
8   try {
9     Path path = new File(LOG_FILE).toPath();
10    String content = Files.readString(path);
11    return content;
12  } catch (IOException e) {
13    throw new RuntimeException(e.getMessage());
14  }
15 }

```

Fig. 12. Example method with a bug on Line 3 which should instead be `if(ready)`. The program crashes upon executing Line 13.

frames. As a result, PROSECUTOR ranks the bugs in versions #16 and #41 of the Compress project [2] at Rank 1, whereas SmartFL places them at ranks 62 and 358, respectively.

Example 4. Finally, our constructed EPG also models the effect of unexplored branches on the executed statements. Consider the example code in Figure 13 where given two input strings that represent two integers with different number of digits, the method `getMinStr` produces the integer value for the number with fewer digits. Given the failing execution in this figure, our inference rules derive and include the following edges in the EPG: $\langle \# @ 8 \rangle \rightarrow \langle \text{str} @ 9 \rangle$, $\langle \text{val} @ 11 \rangle \rightarrow \langle \text{res} @ 9 \rangle$, and $\langle \# @ 8 \rangle \rightarrow \langle \text{res} @ 9 \rangle$. Similar to the previous example, we also incorporate the observations $\neg \langle \text{res} @ 9 \rangle$ and $\neg \langle \text{val} @ 11 \rangle$ into the model which allows it to place the `if`-statement on Line 10 within its top predictions for the bug location.

On the other hand, none of these edges are included in SmartFL’s probabilistic graph since (a) SmartFL does not consider the effect of unexplored branches which we handle by Rules R_{OP3} , R_{OP4} and R_{CTR2} , and (b) it does not model the affect of callee statements on the caller statements when the callee procedure terminates ungracefully due to its reliance on actual stack execution.

This example demonstrates how PROSECUTOR achieves more accurate bug localization (as a result of (a)) in Jsoup version #39, placing the actual faulty statement at Rank 7, while SmartFL places it at Rank 632. At the same time, it demonstrates how considering (b) enables PROSECUTOR to localize the fault in version #1 of the Lang project at Rank 2 while SmartFL places this bug at Rank 8 of its ranking.

Example 5. In contrast to the example in Figure 11, note that the set of shortcut edges we draw in the EPG can occasionally have a negative effect on the fault ranking. Consider the code in Figure 14 with a bug in the method `absURL`. In this example, our inference rules include the following two edges in the EPG: $\langle a1 @ 2 \rangle \rightarrow \langle \text{str} @ 3 \rangle$ and $\langle \# @ 6 \rangle \rightarrow \langle \text{str} @ 3 \rangle$. Notice that SmartFL’s graph does not include the edge $\langle a1 @ 2 \rangle \rightarrow \langle \text{str} @ 3 \rangle$. Since the actual fault resides on the path reachable by tracing backward through the edge $\langle \# @ 6 \rangle \rightarrow \langle \text{str} @ 3 \rangle$, SmartFL ranks the fault at Rank 2, while Prosecutor spreads suspicion somewhat between the two paths and—after Bayesian inference—ranks the actual fault at Rank 7. Removing shortcut edges in this case can raise the bug’s position in the ranking to Rank 4 instead.

```

1 public static Integer getMinStr(String a, String b) {
2   if (a == null || b == null)
3     throw new IllegalArgumentException("Null Args");
4   int a_len = a.length();
5   int b_len = b.length();
6   int min = Math.min(a_len, b_len);
7   if (min > 10 || a_len == b_len)
8     throw new IllegalArgumentException("Invalid Args");
9   String s = a;
10  if (b_len > a_len) // faulty line
11    s = b;
12  Integer res = createInteger(s);
13  return res;
14 }
15 public static Integer createInteger(final String str) {
16   if (str == null)
17     return null;
18   Integer val = Integer.decode(str); // crash
19   return val;
20 }
21
22 public void testCase() {
23   Integer x = StringUtils.getMinStr("12345678912", "12");
24   assertEquals(12, x);
25 }

```

Fig. 13. Example string processing method with a bug on Line 10 which must instead be `if(b_len < a_len)`. This synthesized example shows how PROSECUTOR more effectively localizes bugs in version #39 of the Jsoup project (arising from an unexplored branch) and the one in version #1 of the Lang project (located in a parent method) in Defects4J.

```

1 public void testCase() {
2   Document doc = Jsoup.parse(
3     "<a href=?foo?>One</a> <a href='bar.html?foo?>Two</a>",
4     "http://jsoup.org/path/file?bar");
5   Element a1 = doc.select("a").first();
6   String str = a1.absUrl("href")
7   assertEquals("http://jsoup.org/path/file?foo", str);
8 }
9 public String absUrl(String attributeKey) {
10  ... /* several statements inside the method */
11  URL abs = new URL(base, relUrl);
12  String res = abs.toExternalForm();
13  return res;
14 }

```

Fig. 14. Test snippet adapted from the Jsoup library [4] to describe the bug in version #10 of this project in Defects4J. This code has a bug on Line 11, which results in an incorrect value being returned.

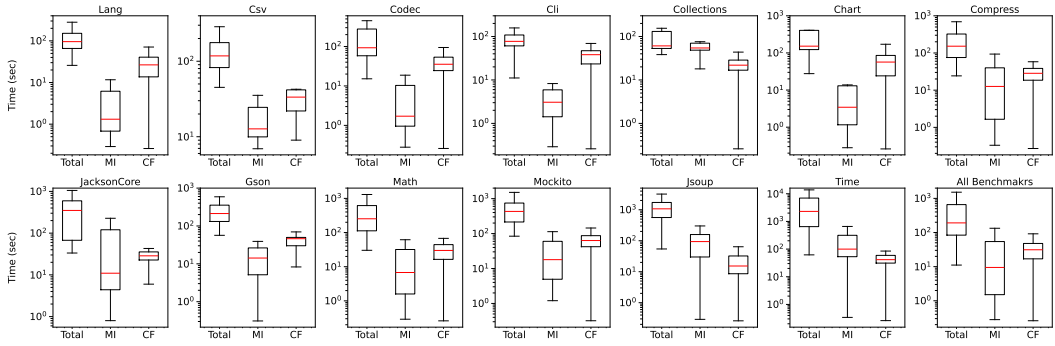


Fig. 15. PROSECUTOR running time breakdown, including the total time and the time needed for performing for marginal inference (MI) and counterfactual analysis (CF).

Overall, note that SmartFL does not involve several types of edges which we consider to make fault localization more general and effective across diverse scenarios. In benchmarks where the bug can be captured only through simple data and control dependencies—where SmartFL’s modeling is already sufficient—the additional edges introduced in the EPG may occasionally have a slight negative impact on the ranking results. However, as we demonstrate in our ablation study, these edges significantly improve the generality of PROSECUTOR. With all edges included, PROSECUTOR localizes more than 73% of the faults within the top 20 entries of its ranking, while removing shortcut edges at return sites and call sites reduces this number to 66% and 69%, respectively.

5.2 RQ2: Runtime Overhead of PROSECUTOR

We also measured the total time that each approach needs to provide its ranking of possible fault locations. PROSECUTOR requires an average of 240 seconds to localize each fault (using geometric means). EPG construction (56%), marginal inference (15%) and counterfactual analysis (5%) turn out to be the most significant contributors to this running time. Figure 15 shows the breakdown of these statistics across different projects.

Unsurprisingly, we observed that SBFL techniques are consistently the most lightweight, and MBFL approaches are the most expensive, requiring on average 75 and 770 seconds per benchmark, respectively. We also observed that SmartFL is $\approx 3\times$ faster than our approach with a speed comparable to SBFL. Most of this gap in running time is due to less accurate and faster graph construction, for which SmartFL provides a parallelized implementation. In principle, we can also parallelize EPG construction, which we expect to lead to similar speedups.

Our implementation for PROSECUTOR provides an average graph construction throughput of 70 stmt/sec, and the loop identification algorithm could reduce the length of traces with an average speed of 710 stmt/sec. Furthermore, our experiments show that the loop reduction on average eliminates more than 38,000 statements from the processed traces in each benchmark. With loop compression turned off, PROSECUTOR requires an average of 500 seconds per benchmark to produce its ranking, which is $\approx 2\times$ the time needed by PROSECUTOR with loop compression turned on. This clearly highlights the value of loop identification in the overall scalability of the procedure.

5.3 RQ3: Ablation Study

5.3.1 Effect of Passing and Failing Traces. We ran PROSECUTOR in two ablated settings to measure the impact of failing and passing traces on the eventual ranking quality. We thus eliminate the ranking aggregation phase and individually consider the two rankings obtained by constructing EPG from failing traces, F , and by using both failing and passing traces, $F \cup P$.

We aggregated these statistics using the top- k measure in Table 5. By comparing the statistics reported in Tables 3 and 5, one can conclude that even the first setting, which does not leverage passing traces, outperforms the best-performing baseline in the top k predictions for $k \in \{1, 3, 5, 10, 20\}$. Note that in Defects4J, on average less than 2 failing tests trigger the bug for each benchmark. Thus, one notable advantage of our approach is its ability to effectively localize the fault even without access to passing traces by using a few failing tests. This is essentially impossible when using SBFL techniques. Observe also that by incorporating the dependency information from passing traces in the EPG, PROSECUTOR can further enhance the quality of its ranking.

5.3.2 Effect of Loop Compression. Since our loop compression algorithm eliminates redundant statements (and therefore redundant dependencies) from execution traces, we also conducted an ablation study to assess the possible impact of this compression on the overall accuracy of Prosecutor. As shown in Table 5, the results of our study indicate that loop compression not only does not negatively impact the accuracy of PROSECUTOR but also slightly improves it, while still significantly enhancing the system’s efficiency.

5.3.3 Effect of Shortcut Edges. We also conducted two additional ablation studies where we disabled the shortcut edges described in Section 5.1.2 to assess their impact on ranking accuracy.

Our first experiment corresponds to disabling Rules R_{CL1} when t_i is not a **return** statement. In this setting, we also exclude the set of observations on the incorrectness of $\{v \in V_i^d \mid \text{CLOSE}_i \wedge \text{OPEN}_i\}$, where v represent destination/target vertices in these shortcut edges.

Our second experiment involves changing the Rule R_{DU} to the following rule, which simply exclude shortcut edges that abstractly account for the potential impact of method arguments and the base object on the return value at call sites:

$$R'_{DU} \frac{v \in V_i^d \quad u \in \text{USE}(l_i) \quad \neg \text{ISRCVObj}_i(u) \quad \neg \text{ISARG}_i(u) \quad \langle u@j \rangle = \text{LASTDEF}(u, t_i)}{\langle u@j \rangle \rightarrow v \in E}$$

Once again, Table 5 presents the results of this experiment. Overall, observe that removing shortcut edges in our first experiment decreases the top-1, top-3, and top-5 metrics consistently by 16%. As described earlier, these edges are specifically derived by R_{CL1} when the execution fails within a nested callee procedures, where several call frames are alive within the execution stack (the examples in Figure 1, 12, and 13 are among such cases). These failures, which mainly occur in the form of unexpected exceptions or program crashes, leads to abrupt execution termination. Our measurements show that, among the 867 failure-triggering tests in our benchmark set, 423 (49%) tests fail before reaching assertions in the corresponding test procedure. Since these cases account for a large fraction of the failing tests, disabling shortcut edges at return sites naturally results in a noticeable decline in accuracy.

Observe also that removing shortcut edges at invocation sites has a limited impact on the top-10 metric, yet it still significantly affects the top-20 metric.

5.3.4 Effect of Control Dependency Edges. As our final ablation study, we excluded the control dependency edges from the EPG to evaluate how data and control dependencies contribute to fault localization performance. We specifically disabled the Rule R_{CTR1} , R_{CTR2} , R_{OP3} , and R_{OP4} when

Table 5. Effectiveness of PROSECUTOR in various ablated settings. The rows labeled F and $F \cup P$ show the performance without ranking aggregation, and the subsequent rows present the results without loop compression, shortcut edges, control dependency edges and counterfactual analysis, respectively.

Setting	TOP-1	TOP-3	TOP-5	TOP-10	TOP-20	EXAM
SmartFL	95	172	201	236	271	10.2
PROSECUTOR	106	189	231	289	342	4.2
F	98	182	220	272	323	4.7
$F \cup P$	103	189	232	280	326	4.7
w/o Loop Compression	103	185	227	286	339	4.4
w/o Shortcut (Call site)	104	183	229	284	324	4.8
w/o Shortcut (Return site)	89	159	196	256	310	6.0
w/o Ctrl Dependencies	85	167	202	262	302	5.5
w/o Counterfactual	104	183	221	279	332	4.5

1079 constructing the EPG from execution traces. Notice also that in this study, we have still preserved
 1080 value vertices which represent the correctness of the predicates evaluated along the traces. We
 1081 aggregated the results under this setting and report them in Table 5.

1082 Observe that without control dependency edges, PROSECUTOR localizes 20% and 12% fewer buggy
 1083 lines within its top 1 and top 5 predictions, respectively. Note also that excluding these edges may
 1084 naturally improve the rankings for bugs that only affect data along a failing execution. This is
 1085 because disabling control dependency edges eliminates certain paths between vertices in the EPG,
 1086 thereby potentially reducing the set of statements reachable from the failure location and, in turn,
 1087 shrinking the search space of candidates for the actual bug location.

1088

1089 5.4 RQ4: Role of Counterfactual Executions

1090 To assess the impact of counterfactual analysis on the overall effectiveness of the rankings, we
 1091 performed a final experiment with the feature turned off. Table 5 also includes the results of this
 1092 study. Observe the number of faults identified within the top- k entries of the ranking drops without
 1093 counterfactual observations.

1094 In particular, we observed that in 38% of the benchmarks, at least one of our counterfactual
 1095 experiments change the eventual test outcome. This clearly shows that the initial Bayesian inference
 1096 can effectively prioritize suspicious branch conditions potentially related to the fault, which again
 1097 highlights the effectiveness and precision of our modeling. In such cases, on average, counterfactual
 1098 executions improved bug rank by 4 positions. 77% of bugs were positively affected or unchanged.
 1099 Bugs that were positively affected saw their ranks improve by 25 positions, on average, while
 1100 bugs that were negatively affected only saw their ranks drop by 7 positions. Furthermore, the
 1101 successful counterfactual experiments also increased the number of faults identified in top 3 and
 1102 top 5 predictions by 7% and 10%, respectively.

1103 We specifically observed that counterfactual analysis
 1104 helps in raising the priority of deep bugs in the ranking.
 1105 While the median distance of buggy statement vertices
 1106 from the failure location in the EPG is 6 edges, it rises to
 1107 9 in cases where counterfactual analysis improves the
 1108 ranking. From a high-level view, program elements that
 1109 are far away in the EPG from the failure point typically
 1110 appear less suspicious. On the other hand, marking
 1111 the virtual variable `cfx` as incorrect draws suspicion
 1112 towards the statement nodes with edges leading to `cfx`
 1113 vertex and their surrounding statements.

1114 For example, version #17 of the Cli project in De-
 1115 fects4J has a bug that is 24 edges away from assertion
 1116 violation in the EPG. This causes the actual fault to
 1117 be placed in Rank 52 in the ablated setting, while the
 1118 full PROSECUTOR places this statement at Rank 16. This
 1119 evidently shows the effect of counterfactual analysis in
 1120 prioritizing deep bugs.

1121 A broader view of this improvement across all benchmarks with deep bugs is illustrated in
 1122 Figure 16. We define bug depth as the number of edges in the EPG from the buggy statement to the
 1123 failure location. We consider bugs whose depths exceed the third quartile (i.e., ≥ 11 edges) as “*deep*
 1124 bugs”. Compared to SmartFL, PROSECUTOR without counterfactual analysis can localize 50% more
 1125 deep faults within its top 10 predictions. Incorporating counterfactual analysis into the method
 1126

1127

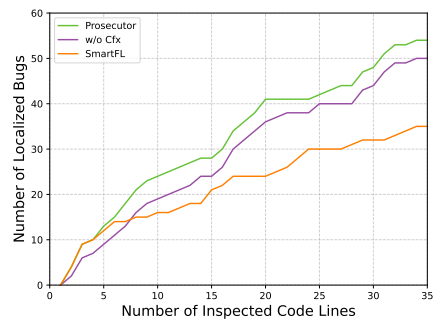


Fig. 16. Number of deep faults in top- k predictions of PROSECUTOR—with and without counterfactual analysis—against SmartFL. The benchmarks totally contain 106 deep faults. The faster rising curves are better.

1128 further enhances the accuracy of localization and allows our approach to identify 50% and 71%
1129 more deep bugs than SmartFL within its top 10 and 20 predictions, respectively.

1130 Finally, beyond improvements in the top- k metric, our experience suggests that counterfactual
1131 analysis serves as an additional and meaningful source of validation. Even in benchmarks where the
1132 true fault location was already ranked within the top predictions prior to applying counterfactual
1133 reasoning, a successful counterfactual run further increased the suspiciousness scores of the faulty
1134 statements, reinforcing confidence in their ranking. We notably, observed that incorporating
1135 counterfactual results back into the model alters the suspiciousness of buggy statements on average
1136 $5\times$ more than that of non-buggy statements.

1137

1138 6 Limitations of Our Approach

1139 We now briefly discuss some limitations of PROSECUTOR that suggest directions for future research.

1140 First, our approach is targeted towards “*errors of commission*”, e.g., errors involving incorrect
1141 values, assertion violations, program crashes, and unexpected exceptions. In contrast, some bugs
1142 manifest as “*errors of omission*”, mainly in the form of expected exceptions that not actually raised.
1143 We handle such cases by performing a lightweight analysis to identify and annotate possible failure
1144 points where the exception might arise. A general solution to this problem is a good direction for
1145 future work.

1146 Next, recall that the EPG provides a model for how errors arise during program executions.
1147 Although the EPG successfully captures a large number of errors occurring within the Defects4J
1148 programs, given that our inference rules are not sound, it is unable to model complex patterns of
1149 aliasing (our current solution uses static specifications to identify possibly aliasing statements)
1150 or adequately model the heap. More precise modeling of data dependencies would presumably
1151 improve the overall effectiveness. It is likely, however, that soundly capturing more sources of
1152 error propagation would, in a manner similar to precision-recall tradeoffs in program analysis and
1153 machine learning, lead to reduced ranking quality overall.

1154 Finally, our current implementation of PROSECUTOR faces difficulties in scaling to traces with
1155 millions of statements. In the Defects4J dataset, this mostly arises from programs with repeated
1156 irregular invocations of the same iterative computation, which limits the effectiveness of the trace
1157 compression procedure. Note that the EPG construction algorithm requires $O(n)$ time and makes a
1158 single pass over the execution trace with n statements. The main performance bottleneck in our
1159 implementation lies in the I/O overhead of exchanging data between Java and Python processes.
1160 Besides parallelization of the graph construction, performing static dependency analysis beforehand
1161 in an offline process can significantly reduce I/O overhead and further enhance the scalability.
1162 Another interesting direction of future work is to identify and prune unnecessary dependencies in
1163 the EPG to make the fault localization both faster and more effective.

1164

1165 7 Related Work

1166 There is a large body of research on fault localization and probabilistic program reasoning. We now
1167 summarize the most closely related work. For a detailed survey, we refer readers to [33].

1168

1169 *Fault localization.* The introduction of delta debugging for simplifying and isolating failure-
1170 inducing inputs [37] was perhaps among the earliest work that sparked a line of research towards
1171 both automated debugging and fault localization techniques. Research on fault localization began
1172 with spectrum-based approaches, which rank program statements by comparing their execution
1173 frequencies in passing and failing tests [6, 7, 13, 25, 32]. It was later developed by mutation-based
1174 techniques, which inspect test behaviors across a set of mutants to assess the impact of changing
1175 each statement on test outcomes [10, 24, 26]. Predicate switching [40] alters control flow by

1176

1177 heuristically identifying critical predicates and inverting conditional outcomes during failing runs,
1178 which may be regarded as a simple form of counterfactual analysis. In contrast, our work describes
1179 how to effectively incorporate the results of these experiments into a single unified probabilistic
1180 model. History-based techniques prioritize fault candidates based on historical defect data [15, 29]
1181 instead of relying on information from test runs. This data may be combined with coverage profiles
1182 to make the fault identification more accurate [31].

1183 Besides statistical approaches, learning-based fault localization has recently gained significant
1184 attention. UniVal [18] uses statistical causal inference to estimate the average causal effect of
1185 counterfactual assignments to program variables, without actually executing the program. Neural-
1186 MBFL [12] uses a pre-trained model that can leverage the context information surrounding the
1187 mutation position to improve the quality of mutants generated. GMBFL [34] models the relationship
1188 between code entities, mutants, and test runs in a graph representation and utilizes gated graph
1189 attention neural networks to learn useful features from this graph.

1190 Another line of work on learning-based method-level fault localization was started by DeepFL [20],
1191 which was further enhanced by newer classifiers. DeepFL incorporates suspiciousness scores of
1192 SBFL and MBFL with code complexity and text similarity into a unified deep-learning model.
1193 Grace [21] leverages graph-based representation learning to fully utilize coverage information that
1194 involves connective relationships between tests and program entities. DepGraph [27] proposes a
1195 more compact representation by integrating the inter-procedural call graph into the GNN model
1196 alongside code change information as an additional feature. HetFL [8] mitigates data imbalance
1197 and inefficient representation of previous methods by resampling faulty data and modeling pro-
1198 grams as heterogeneous graphs. To our knowledge, most of these techniques perform method-level
1199 localization. This limits their effectiveness when applied to failing unit tests, which typically only
1200 execute a small number of functions (In our experiments, failing tests in Defects4j dataset executed
1201 a median of 19 and an average of 30 functions before crashing). In addition, the non-reliance on
1202 training data and the non-necessity of a GPU makes our technique much more broadly applicable
1203 and less expensive than most learning-based approaches.

1204
1205 *Probabilistic program reasoning.* Starting with Eugene [22] and Bingo [28], there is a growing body
1206 of work on using Bayesian inference to probabilistically reason about programs. Xu et al. model
1207 debugging as a probabilistic inference problem, where human-like reasoning rules are modeled as
1208 conditional probability distributions [36]. Inspired by this work, SmartFL [38] constructs an efficient
1209 probabilistic model of program semantics using dynamic data and static control dependencies,
1210 which presents a promising direction, but the proposed encoding of program semantics is subject
1211 to several limitations as mentioned throughout this paper.

1212

1213 8 Conclusion

1214 In this paper, we presented PROSECUTOR, a new fault localization technique that combines Bayesian
1215 reasoning about erroneous traces with a novel counterfactual analysis to effectively identify faults
1216 at the statement level. A comprehensive evaluation on benchmarks from the Defects4J dataset of
1217 faulty Java libraries indicates that our technique is applicable to real-world programs and sets up a
1218 new state-of-the-art in ranking effectiveness.

1219

1220 Data Availability

1221 We have submitted our replication package as supplementary material for the review process.
1222 We will also submit the artifact for artifact evaluation and make it publicly available upon paper
1223 acceptance.

1224

1225

References

- [1] [n. d.]. *Apache Commons Collections java library*. <https://commons.apache.org/proper/commons-collections>
- [2] [n. d.]. *Apache Commons Compress java library*. <https://commons.apache.org/proper/commons-compress>
- [3] [n. d.]. *Apache Commons Lang java library*. <https://commons.apache.org/proper/commons-lang>
- [4] [n. d.]. *jsoup: Java HTML Parser*. <https://jsoup.org>
- [5] [n. d.]. *The Major Mutation Framework*. <https://mutation-testing.org>
- [6] Rui Abreu, Peter Zoetewij, Rob Golsteijn, and Arjan van Gemund. 2009. A Practical Evaluation of Spectrum-Based Fault Localization. *Journal of Systems and Software* 82, 11 (Nov. 2009), 1780–1792. doi:10.1016/j.jss.2009.06.035
- [7] Rui Abreu, Peter Zoetewij, and Arjan van Gemund. 2009. Spectrum-Based Multiple Fault Localization. In *IEEE/ACM International Conference on Automated Software Engineering (ASE)*. 88–99. doi:10.1109/ASE.2009.25
- [8] Xin Chen, Tian Sun, Dongling Zhuang, Dongjin Yu, He Jiang, Zhide Zhou, and Sicheng Li. 2024. HetFL: Heterogeneous Graph-Based Software Fault Localization. *IEEE Transactions on Software Engineering* 50, 11 (2024), 2884–2905. doi:10.1109/TSE.2024.3454605
- [9] Jong-Deok Choi and Jeanne Ferrante. 1994. Static Slicing in the Presence of Goto Statements. *ACM Transactions on Programming Languages and Systems* 16, 4 (July 1994), 1097–1113. doi:10.1145/183432.183438
- [10] Zhanqi Cui, Minghua Jia, Xiang Chen, Liwei Zheng, and Xiulei Liu. 2020. Improving Software Fault Localization by Combining Spectrum and Mutation. *IEEE Access* 8 (2020), 172296–172307. doi:10.1109/ACCESS.2020.3025460
- [11] Ron Cytron, Jeanne Ferrante, Barry Rosen, Mark Wegman, and Kenneth Zadeck. 1991. Efficiently Computing Static Single Assignment Form and the Control Dependence Graph. *ACM Transactions on Programming Languages and Systems* 13, 4 (Oct. 1991), 451–490. doi:10.1145/115372.115320
- [12] Bin Du, Baolong Han, Hengyuan Liu, Zexing Chang, Yong Liu, and Xiang Chen. 2024. Neural-MBFL: Improving Mutation-Based Fault Localization by Neural Mutation. In *2024 IEEE 48th Annual Computers, Software, and Applications Conference (COMPSAC)*. 1274–1283. doi:10.1109/COMPSAC61105.2024.00168
- [13] James Jones and Mary Jean Harrold. 2005. Empirical Evaluation of the Tarantula Automatic Fault-Localization Technique. In *Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. ACM, 273–282. doi:10.1145/1101908.1101949
- [14] René Just, Darious Jalali, and Michael Ernst. 2014. Defects4J: A Database of Existing Faults to Enable Controlled Testing Studies for Java Programs. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA)*. ACM, 437–440. doi:10.1145/2610384.2628055
- [15] Sunghun Kim, Thomas Zimmermann, E. James Whitehead Jr., and Andreas Zeller. 2007. Predicting Faults from Cached History. In *29th International Conference on Software Engineering (ICSE)*. 489–498. doi:10.1109/ICSE.2007.66
- [16] Daphne Koller and Nir Friedman. 2009. *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press.
- [17] Daphne Koller and Nir Friedman. 2009. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press.
- [18] Yiğit Küçük, Tim Henderson, and Andy Podgurski. 2021. Improving Fault Localization by Integrating Value and Predicate Based Causal Inference Techniques. In *Proceedings of the 43rd International Conference on Software Engineering (ICSE)*. IEEE Press, 649–660. doi:10.1109/ICSE43902.2021.00066
- [19] Tom Leonard, John S. J. Hsu, and Kam-Wah Tsui. 1989. Bayesian Marginal Inference. *J. Amer. Statist. Assoc.* 84, 408 (1989), 1051–1058. arXiv:<https://www.tandfonline.com/doi/pdf/10.1080/01621459.1989.10478871> doi:10.1080/01621459.1989.10478871
- [20] Xia Li, Wei Li, Yuqun Zhang, and Lingming Zhang. 2019. DeepFL: integrating multiple fault diagnosis dimensions for deep fault localization. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis (Beijing, China) (ISSTA 2019)*. Association for Computing Machinery, New York, NY, USA, 169–180. doi:10.1145/3293882.3330574
- [21] Yiling Lou, Qihao Zhu, Jinhao Dong, Xia Li, Zeyu Sun, Dan Hao, Lu Zhang, and Lingming Zhang. 2021. Boosting coverage-based fault localization via graph-based representation learning. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Athens, Greece) (ESEC/FSE 2021)*. Association for Computing Machinery, New York, NY, USA, 664–676. doi:10.1145/3468264.3468580
- [22] Ravi Mangal, Xin Zhang, Aditya Nori, and Mayur Naik. 2015. A User-Guided Approach to Program Analysis. In *Proceedings of the 10th Joint Meeting on Foundations of Software Engineering (ESEC/FSE)*. ACM, 462–473.
- [23] Joris Mooij. 2010. libDAI: A Free and Open Source C++ Library for Discrete Approximate Inference in Graphical Models. *Journal of Machine Learning Research* 11, 74 (2010), 2169–2173. <http://jmlr.org/papers/v11/mooij10a.html>
- [24] Seokhyeon Moon, Yunho Kim, Moonzo Kim, and Shin Yoo. 2014. Ask the Mutants: Mutating Faulty Programs for Fault Localization. In *IEEE Seventh International Conference on Software Testing, Verification and Validation*. 153–162. doi:10.1109/ICST.2014.28

- 1275 [25] Lee Naish, Hua Jia Lee, and Kotagiri Ramamohanarao. 2011. A Model for Spectra-Based Software Diagnosis. *ACM*
1276 *Transactions on Software Engineering and Methodology* 20, 3, Article 11 (Aug. 2011), 32 pages. doi:10.1145/2000791.
1277 2000795
- 1278 [26] Mike Papadakis and Yves Le Traon. 2015. Metallaxis-FL: Mutation-Based Fault Localization. *Journal of Software:*
1279 *Testing, Verification and Reliability* 25, 5–7 (Aug. 2015), 605–628. doi:10.1002/stvr.1509
- 1280 [27] Md Nakhla Rafi, Dong Jae Kim, An Ran Chen, Tse-Hsun (Peter) Chen, and Shaowei Wang. 2024. Towards Better
1281 Graph Neural Network-Based Fault Localization through Enhanced Code Representation. *Proc. ACM Softw. Eng.* 1,
1282 FSE, Article 86 (July 2024), 23 pages. doi:10.1145/3660793
- 1283 [28] Mukund Raghthaman, Sulkeha Kulkarni, Kihong Heo, and Mayur Naik. 2018. User-Guided Program Reasoning Using
1284 Bayesian Inference. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and*
1285 *Implementation. SIGPLAN Notices* 53, 4, 722–735. doi:10.1145/3296979.3192417
- 1286 [29] Jeongju Sohn and Shin Yoo. 2017. FLUCCS: Using Code and Change Metrics to Improve Fault Localization. In
1287 *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM, 273–283.
1288 doi:10.1145/3092703.3092717
- 1289 [30] Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie Hendren, Patrick Lam, and Vijay Sundaresan. 1999. Soot: A Java
1290 Bytecode Optimization Framework. In *Proceedings of the Conference of the Centre for Advanced Studies on Collaborative*
1291 *Research (CASCON)*. IBM Press, 13.
- 1292 [31] Ming Wen, Junjie Chen, Yongqiang Tian, Rongxin Wu, Dan Hao, Shi Han, and Shing-Chi Cheung. 2021. Historical
1293 Spectrum Based Fault Localization. *IEEE Transactions on Software Engineering* 47, 11 (2021), 2348–2368. doi:10.1109/
1294 TSE.2019.2948158
- 1295 [32] Eric Wong, Vidroha Debroy, Ruizhi Gao, and Yihao Li. 2014. The DStar Method for Effective Software Fault Localization.
1296 *IEEE Transactions on Reliability* 63, 1 (2014), 290–308. doi:10.1109/TR.2013.2285319
- 1297 [33] W. Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, Franz Wotawa, and Dongchen Li. 2023. *Software Fault*
1298 *Localization: An Overview of Research, Techniques, and Tools*. John Wiley & Sons, Ltd, Chapter 1, 1–117.
1299 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781119880929.ch1 doi:10.1002/9781119880929.ch1
- 1300 [34] Shumei Wu, Zheng Li, Yong Liu, Xiang Chen, and Mingyu Li. 2023. GMBFL: Optimizing Mutation-Based Fault
1301 Localization via Graph Representation. In *IEEE International Conference on Software Maintenance and Evolution*
1302 *(ICSM)*. 245–257. doi:10.1109/ICSM58846.2023.00033
- 1303 [35] Yiqian Wu, Yujie Liu, Yi Yin, Muhan Zeng, Zhentao Ye, Xin Zhang, Yingfei Xiong, and Lu Zhang. 2025. Smartfl:
1304 Semantics based probabilistic fault localization. *IEEE Transactions on Software Engineering* (2025).
- 1305 [36] Zhaogui Xu, Shiqing Ma, Xiangyu Zhang, Shuofei Zhu, and Baowen Xu. 2018. Debugging with Intelligence via
1306 Probabilistic Inference. In *40th IEEE/ACM International Conference on Software Engineering (ICSE)*. 1171–1181. doi:10.
1307 1145/3180155.3180237
- 1308 [37] Andreas Zeller and Ralf Hildebrandt. 2002. Simplifying and Isolating Failure-Inducing Input. *IEEE Transactions on*
1309 *Software Engineering* 28, 2 (2002), 183–200. doi:10.1109/32.988498
- 1310 [38] Muhan Zeng, Yiqian Wu, Zhentao Ye, Yingfei Xiong, Xin Zhang, and Lu Zhang. 2022. Fault Localization via Efficient
1311 Probabilistic Modeling of Program Semantics. In *Proceedings of the 44th International Conference on Software Engineering*
1312 *(ICSE)*. ACM, 958–969. doi:10.1145/3510003.3510073
- 1313 [39] Xin Zhang, Radu Grigore, Xujie Si, and Mayur Naik. 2017. Effective interactive resolution of static analysis alarms.
1314 *Proc. ACM Program. Lang.* 1, OOPSLA, Article 57 (Oct. 2017), 30 pages. doi:10.1145/3133881
- 1315 [40] Xiangyu Zhang, Neelam Gupta, and Rajiv Gupta. 2006. Locating Faults Through Automated Predicate Switching. In
1316 *Proceedings of the 28th International Conference on Software Engineering (ICSE)*. ACM, 272–281. doi:10.1145/1134285.
1317 1134324