

Understanding Differences among Executions with Variational Traces

Jens Meinicke,^{1,2} Chu-Pan Wong,¹ Christian Kästner,¹ Gunter Saake²

¹Carnegie Mellon University, USA, ²University of Magdeburg, Germany

Abstract—One of the main challenges of debugging is to understand why the program fails for certain inputs but succeeds for others. This becomes especially difficult if the fault is caused by an interaction of multiple inputs. To debug such interaction faults, it is necessary to understand the individual effect of the input, how these inputs interact and how these interactions cause the fault. The differences between two execution traces can explain why one input behaves differently than the other. We propose to compare execution traces of all input options to derive explanations of the behavior of all options and interactions among them. To make the relevant information stand out, we represent them as variational traces that concisely represents control-flow and data-flow differences among multiple concrete traces. While variational traces can be obtained from brute-force execution of all relevant inputs, we use variational execution to scale the generation of variational traces to the exponential space of possible inputs. We further provide an Eclipse plugin Varviz that enables users to use variational traces for debugging and navigation. In a user study, we show that users of variational traces are more than twice as fast to finish debugging tasks than users of the standard Eclipse debugger. We further show that variational traces can be scaled to programs with many options.

Index Terms—Debugging, Program Comprehension, Feature Interaction, Configurable Software, Variational Execution.



1 INTRODUCTION

Understanding why a certain program input causes a fault while another succeeds is a common task during debugging [89]. This happens, for example, if a certain program crashes in one configuration but succeeds if a parameter is changed [19], [54], [61]. Reasoning about such differences in executions is difficult when using a standard debugger as the program can only be executed for one input at a time. To understand why certain inputs lead to a fault requires to understand the differences between the valid and the invalid executions. Such a comparison can be generated by recording and aligning the traces and state changes of the two executions [80], [87]. The aligned traces can be used to generate explanations of the fault [81], [90].

Some faults, however, are caused by interactions of multiple inputs which make understanding and debugging them even more challenging (aka. *feature interaction*) [1], [3], [15], [22]. Such interaction faults are common in practice, especially in highly configurable systems with many options [1], or if a set of changes introduces the fault [9], [24], [91]. Interaction faults are hard to detect as they require to specify a certain input to trigger the fault [55]. Even if we can narrow down the fault to a smaller number of options, say with *delta debugging* [89], [91], it is still difficult to understand *why*, *where*, and *how* they interact.

After identifying the set of interacting options, a programmer can start investigating how this interaction causes the fault. Understanding the interaction requires understanding the individual behavior of the interacting options, but also their combined behavior. Thus, it may no longer be sufficient to align the execution of two inputs as previous approaches do [81], [90], as such an alignment cannot explain the effects of multiple options.

We propose to align the execution traces of all configurations to explain the effect of multiple options. We introduce *variational trace*, a compact representation of the

trace differences among all executions. In the variational trace, redundant parts are shared and individual parts are annotated with the input they belong to. This focus on differences allows understanding how data and control flow influence the executions and interact, and thus how the different inputs cause a fault. In Figure 1e, we show a variational trace that illustrates how the interaction fault of the example program is caused, as we will explain.

Generation of variational traces challenges scalability as the number of executions and traces that need to be aligned can potentially be exponential to the number of options. A baseline approach would execute all configurations separately and thus can only scale to explain the interactions among few options. More severely is the potential memory consumption as this approach needs to keep all past statements in memory, as it is upfront unclear which statements will differ among executions [41], [65].

We use variational execution (a.k.a. faceted execution) [7], [56], [57], [60] to avoid separately recording and aligning the traces for many inputs. *Variational execution* runs *all* program configurations simultaneously, often efficiently, by sharing redundancies of the executions and values among these configurations. Since the program is executed only once in a shared fashion, alignment is achieved by construction. Additionally, executed statements and program states are tracked as they relate to the original inputs and their interactions. This means that executions and variables can always be linked to specific program inputs. With variational execution we also avoid the memory explosion as it enables us to decide on the fly which statements differ and need to be kept for the variational trace.

To enable developers to interact with variational traces, we developed an interactive Eclipse plug-in called *Varviz* that visualizes variational traces. *Varviz* can be used for understanding faults, but also for other program comprehension tasks that involve understanding differences

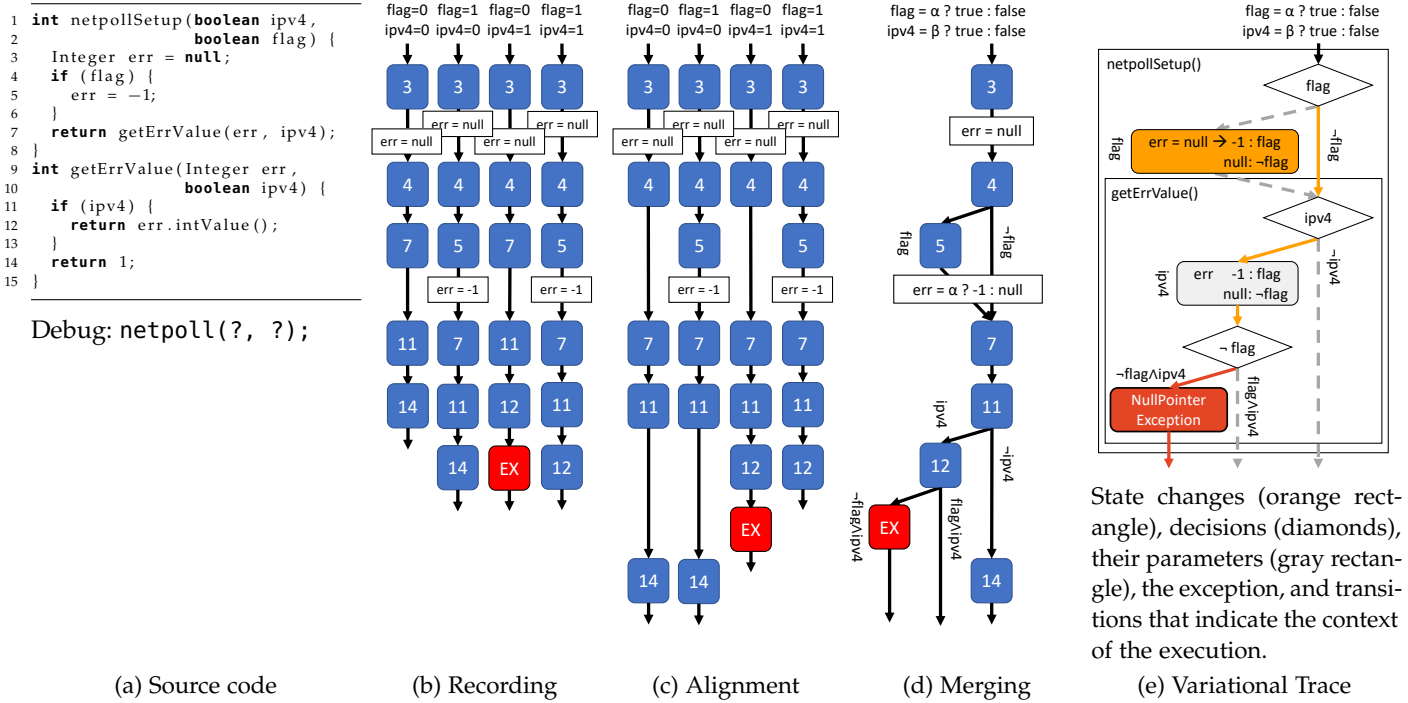


Fig. 1: Concept of a variational trace explaining the interaction fault caused by options `flag` and `ipv4`. Comparing executions requires to *record* the executions for all inputs, *align* the recorded traces, and *merge* the aligned traces to identify interactions (the number on the elements indicate the line number). By enriching the merged trace with information about the state and control-flow decisions, a *variational trace* can be created.

among similar executions.

In a user study, we show that variational traces help to understand a fault in a highly interacting program and that they help to explain a fault from a previous experiment on automatic debugging techniques. Participants using variational traces were more than twice as fast and more successful compared to participants using a standard debugger. In our qualitative study, we found that when dealing with faults caused by differences in inputs, in contrast to standard faults, participants search for places where these inputs interact and how they have an effect, which is exactly what our approach helps with. Furthermore, we show that variational traces are fairly compact, even for medium sized programs with many explored options.

The goal of this research is to aid researchers and developers detecting differences among executions depending on potentially multiple options by providing a dynamic analysis that is able to efficiently execute and align all configurations summarized in variational traces.

Overall, the contributions of this paper are:

- The concept of a *variational trace* that compactly represents differences among the executions of many configurations.
- A *baseline implementation* that shows challenges of generating variational traces.
- *Efficient trace alignment* and a *dynamic analysis to trace only relevant data* avoiding memory explosion for a potentially exponential number of configurations using variational execution [56], [57].
- An Eclipse plug-in *Varviz* to visualize and interact with variational traces.

- A *user study* that shows that participants using variational traces outperform participants using a standard debugger. The study also shows that comparative approaches [81], [90] actually help with debugging tasks.
- An *evaluation on scalability* showing that our approach is able to align executions for large number of input combinations, while concisely describing their differences. By focusing on the effects of certain inputs, the sizes can be reduced even further.

2 STATE OF THE ART

An important problem of many program comprehension and debugging tasks is to understand the differences among executions. For example, a programmer may want to understand why a fault occurs for certain inputs, but does not for others. Such faults can be hard to detect, as they are only triggered for a certain input, and hard to understand because they may require reasoning about interactions among these inputs within long traces. The differences in control flow and data flow among faulty and valid executions can provide insights about how the fault is caused.

Our approach combines ideas from two research fields, namely *automatic debugging* and *feature interactions*, to explain differences among executions. In this section, we discuss how automatic-debugging techniques exploit differences in executions and how various approaches address the feature-interaction problem that causes undesired behavior for certain combinations of inputs.

2.1 Automated Debugging Techniques

Automated-debugging techniques aim to create explanations of why a fault appears, by comparing valid and faulty executions of the program [30], [31], [79], [81], [84], [90]. To create explanations, the approaches execute and record the program multiple times for different inputs or test cases.

Fault Localization. There are many approaches that aim to find the cause of a fault [86]. One of these approaches is *spectrum-based fault localization* which rates each line of the source code by whether it is probable to cause the fault [2], [31]. For example, Tarantula compares the code coverage of valid and failing test cases and provides this information using different background colors on the code [31]. However, such statement ranking has been shown to be less useful than expected [63]: there are usually too many lines highlighted in the code, which makes it hard for users to understand which lines matter for an explanation of the fault. Also, the information why certain lines are important is missing, as are control- and data-flow information. Such information is however often necessary to understand how certain parts of the execution lead to the fault.

Execution Comparison. The comparison of internals of failing and valid executions can be used to explain why programs fail. Instead of just comparing the coverage information, *execution comparison* approaches compare traces and program states across executions [4], [23], [30], [38], [79], [81], [84], [90]. Such comparison can highlight differences in data and control flow relevant to the fault.

Delta debugging is an approach that systematically narrows down the inputs (resp. changes) that are relevant for causing a fault [91]. Based on this idea, Zeller applied delta debugging to program states of executions [90]. A challenge with delta debugging is that it needs to align statements and program states of independent executions [87]. Sumner et al. [79], [81] improved the initial work of Zeller using dual slicing [30], [84] and efficient execution indexing [87]. They improve the efficiency and minimize the explanations in finding the cause effect chain [79], [81].

Execution comparison approaches are designed to explain the differences between only two executions at a time. Thus they require significant runtime overhead, as they execute the program many times to narrow down the instructions necessary for a certain behavior. Due to the separate execution of the program, these approach need to deal with three major challenges, correct alignment of the executions [47], [80], [87], memory overhead that comes with recording the executions [12], [32], [40], and differences in executions caused by nondeterminism [47].

2.2 Understanding Feature Interactions

Feature interaction bugs are hard to detect as they are only triggered for certain combinations of features. There is a lot of research to efficiently find such faults, such as combinatorial interaction testing [19], [54], [61], systematic sampling [34], [37], [76], model checking [10], [18], [83], and variational execution [5], [36], [57], [60].

Sampling approaches can only reveal configurations that fail, but not the interaction that causes it. Variational execution tracks the exact combination of inputs that lead to

a fault, but does not help to understand *why* the interaction happens and *how* it causes the fault.

Other approaches statically reason about the code to detect feature interactions. The work of Kim et al. reasons about the combinatorics to reduce the number of configurations to execute [35]. However, after testing these reduced configuration, the approach cannot answer *why* configurations fail. LoTrack reasons about which lines of code are affected by load-time options [51], but it cannot answer *why* and *how* options interact. The research of Zhang and Ernst aims to identify configuration faults using thin slicing and a lightweight form of execution comparison [77], [92], [93]. Their approach suggests single configuration options that are likely to trigger a fault. However, they assume that (a) the program itself is correct and (b) that a single option triggers the fault rather than a combination of multiple options.

In summary, there exist many approaches that help to compare two executions and approaches that help detecting interactions in larger configuration spaces. However, none of the existing approaches helps to understand *how multiple options interact* as they either do not scale to multiple options or miss control and data-flow information. In our work, we provide support that scales to explain interactions among multiple options and that also provides necessary control and data-flow information about the interactions.

3 GENERATING AND VISUALIZING VARIATIONAL TRACES

In this section, we introduce the new concept of *variational traces* which help developers understanding how options interact during the program execution and assist with debugging interaction faults.

3.1 Variational Traces

We introduce a *variational trace* as a compact representation of differences among multiple execution traces regarding control flow and the program states. With variational traces, we explain how differences in inputs affect a program's execution and data, and how inputs interact with each other. A variational trace is a graph that represents differences on the control-flow and on data using the following concepts (illustrated in Figure 1e for the program in Figure 1a):

- *State changes* (orange rectangles) describe statements that change values of fields and local variables. In the example, the value of *err* is changed from *null* to -1 if the option *flag* is true.
- *Decisions* (diamonds) describe statements causing control-flow differences due to differences in inputs (directly or indirectly) between the individual traces (e.g., if-statements).
- *Decision Parameters* (gray rectangles) describe variables that are used in decisions. In the example, *err* is used as parameter for the decision of Line 12. We found this information on parameters used in decisions particularly useful as they help to understand causes of control-flow differences.
- *Exceptions* (red rectangles) describe statements that are thrown due to faults or exception handling. In the figure, we see that the *NullPointerException* is thrown under the condition that *flag* is false and *ipv4* is true.
- *Return* statements (not shown in the example) describe values that are returned by a method. Return statements

are included if a method returns different values or from different locations.

- *Methods* (rectangles around subgraphs) structure the variational trace and describe the stack trace (e.g., method `getErrValue`). The notation of methods helps to understand the control flow and the execution of the program across method boundaries.

Variational traces show differences. To be a useful debugging tool, variational traces need to concisely describe differences among executions while still containing sufficient information. A variational trace only describes state changes, decisions, invocations and exceptions that differ among executions. Thus, a variational trace contains only statements that cause control-flow differences and statements that change variables differently in different executions. Statements that do not cause such differences are not as important and will not be contained in the variational trace. For example, the statement in Line 3 of Figure 1a sets the value of `err` to `null` in all configurations. This state change is thus not relevant to describe differences among executions. The same applies to changes that only take place on variables that only exist for certain inputs. Thus, if an object only exists under condition *a*, then changes that happen to this object under condition *a* are not important for comparing traces. Instead, the variational trace will report changes that create interactions, such as Line 5 of the example that changes `err` to be either `-1` or `null` depending on the value of `flag`. Such a statement reduction requires to record and keep all state changes of all variables and to compare their values across all executions.

A focus on differences among executions allows to narrow down the number of statements that are reported to the user. Our approach might be further combined with slicing to remove statements that do not matter to explain a certain exception [79], [81], [90]. In this work, we focus on efficiently generating explanations for differences among many executions. Thus, improvements such as data-flow and impact analyses are out of this paper’s scope.

A variational trace, as in Figure 1e, looks similar to a control flow graph. However, three differences make variational traces more useful for debugging purposes: (i) variational traces represent actual executions and they show the actual values that variables take plus the exact context they are executed in, (ii) variational traces only contain statements that differ among the executions, and (iii) variational traces support (beyond others) loops, recursion, dynamic invocation, and reflection.

Using Variational Traces. Variational traces help understanding differences among executions as they describe the cause and effects of differences among executions. This allows answering questions about interactions among options and to understand faults caused by such interactions. First, a variational trace answers for which inputs the interactions occurs. This information is visible in the context of the statement (resp. exception). Second, it helps understanding why the interaction occurs as it shows the differences in the state and in the control-flow at the point in the execution before the statement. Third, a variational trace helps understanding how the interaction is caused and how the corresponding state is created. As an interaction is always caused due to previous control-flow decisions

or other interactions, the information about the cause will always be contained in the variational trace. For example, the exception in Figure 1 is thrown because the value of `err` is null if `flag` is false. The cause of the fault is that `err` was not initialized if `flag` is false, which is visible in the trace.

3.2 Generating Variational Traces by Aligning Trace Logs

To generate variational traces we could align the traces from executing the traces from executing all configurations separately. However, this way of creating variational traces would be challenging. In this subsection, we use such a base line approach that aligns trace logs to explain the necessary steps and challenges of creating variational traces. We illustrate each step with a running example in Figure 1, resulting in the variational trace shown in Figure 1e.

The first step for creating a variational trace is to *record* the program’s executions for *all* different inputs. The recording needs to track both the instructions and the state changes, to understand the effects of the executions. The recordings including state changes for the four inputs of the example are illustrated in Figure 1b.

To compare the executions, all traces need to be aligned after recording. To *align* the traces, instructions from each trace that belong to each other need to be identified [59]. In general, there are multiple different alignments possible as statements may repeat in different parts of the execution. Thus, a semantically optimal alignment is required for meaningful results. However, as the program is executed separately multiple times, alignment is a non-trivial task, especially as object references are different among executions, as it is necessary to keep track of iterations, and as nondeterminism may lead to wrong alignments, which can be avoided by recording and replaying nondeterministic result (e.g., IO) among the executions [47], [80], [87]. We illustrate the alignment of the traces in Figure 1c.

After aligning the traces, shared parts of the executions can be identified. Thus, we can *merge* the traces into a single trace with conditional statements [68]. Statements and state changes that are included in several traces can be shared. We illustrate merging of the four executions from the example in Figure 1d. As shown, Line 5 is only executed if the option `flag` is true, and the value of `err` depends on the input `flag`.

Finally, the variational trace can be *reduced* to statements that describe differences and can be enriched with information about the effects of the instructions. We illustrate the variational trace for the example program in Figure 1e. As shown, the variational trace only describes the control-flow and state differences among the executions resulting in a concise explanation of the differences and the fault.

Discussion. The baseline approach has the obvious issue that it can only scale to few options as it needs to run the program 2^n times for n boolean options. Furthermore, as it is unclear which statements will interact upfront, the approach needs to record full traces which causes severe memory problems. That is, the approach needs to keep all statements and states in memory until all configurations are executed. There are several ways to reduce the memory consumption, such as (i) directly merging the traces after execution instead of keeping them all in memory, (ii) statically deciding which

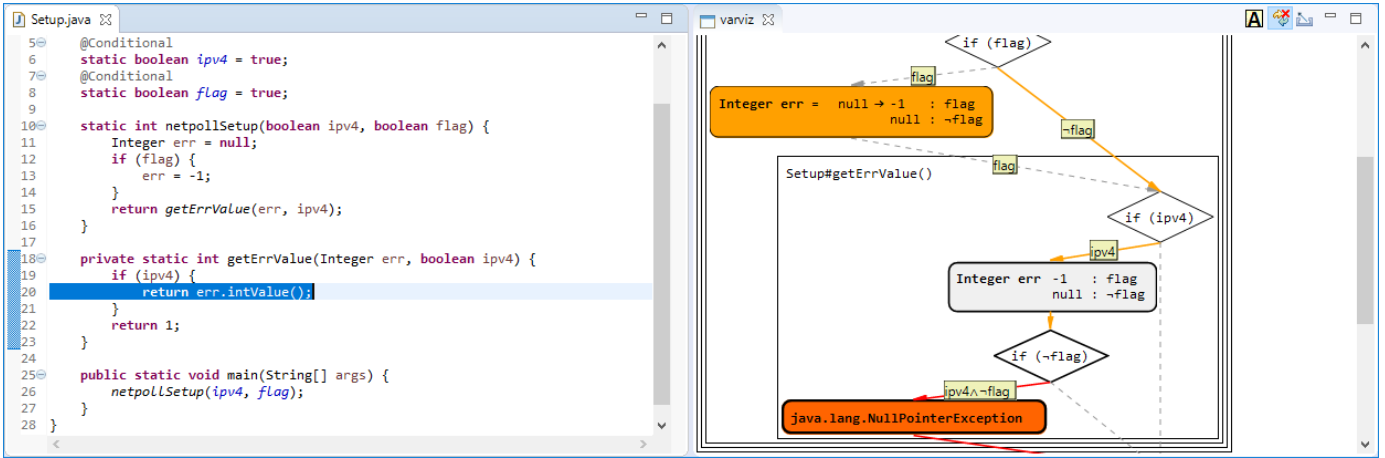


Fig. 2: Screenshot of Varviz in Eclipse for the example of Figure 1a: Java editor (left), variational trace (right).

parts of a method will be equal independent of options (e.g., initialization of local variables), (iii) compressed storage of trace logs [12], [32], and (iv) inter-procedural analysis to detect statements that depend on options [51]. However, these approaches are either not sufficient (i, ii & iii) or do not scale to larger programs (iv).

We implemented the baseline approach in a prototype for Java including optimizations (i) and (ii). We observed that the additional logging statements and the merging process cause a lot of computational overhead. However, especially the memory consumption of the approach is problematic. If too many statements and states must be kept in memory the tool may record multiple GB for the systems in our evaluation, as we show in Table 5, which makes this approach infeasible for larger systems. To approach these challenges requires sophisticated mechanisms to reduce the number of statements to record and to solve the problems of nondeterminism (e.g., by synchronously executing all configurations [38], [47]).

Our goal is to generate variational traces anyhow as they enable us to observe effects and interactions of options. To this end, we have to avoid all the limitations (scalability, memory overhead, nondeterminism) of the base line approach while keeping the idea of aligning the executions of all configurations.

3.3 Efficient Generation of Variational Traces with Variational Execution

A key insight of our work is to use *variational execution* to sidestep the previously discussed challenges for generating variational traces (i.e., handling the challenges of recording, alignment, merging and nondeterminism). Variational execution is an approach that can efficiently execute *all* configurations of the program in a single run [57], [60]. Scalability is achieved because all redundant parts among the executions are executed only once and redundancies in data are stored only once [57]. If parts of the program are only executed under certain conditions, then also variational execution will run these parts only under these conditions. For example, Line 5 of Figure 1a is only executed if the option *flag* is true. Variational execution runs the program with all inputs as illustrated in Figure 1d. Instead of executing parts multiple times as in Figure 1b, each instruction is only executed once.

To represent state differences between configurations, variational execution uses choice values [21]. A choice value is a mapping from partial configurations spaces to concrete values. For example, the value of *err* in Figure 1e is stored as a choice that represents the state of *err* for all four configurations. Using a variational data representation allows variational execution to efficiently represent the states of all configurations as most data is redundant among configurations [56], [57]. Sharing of redundancies in executions and data allows variational execution to scale for programs with huge configuration spaces that would not be practical with a brute-force approach that runs each input individually [57], [60]. We solve the problem of non-deterministic behavior among the configurations due to the sharing and synchronous execution of variational execution, for which an alignment-based approach requires sophisticated synchronization and alignment strategies [47].

Generating Variational Traces. As variational execution synchronizes the executions among configurations, there is no need to align the executions of the program. Instead, we directly generate the variational trace by recording how variational execution runs the program. The context of each statement is already given as all instructions are executed under a certain context. To observe where the execution splits, we just need to observe changes in the execution context. Finally, we can observe interactions on data in the assigned choice values of local variables and fields. With variational execution we avoid the memory explosion of the base line approach, as we can decide during execution whether a statement should be contained in the variational trace. As we will show as part of our evaluation in Figure 5, using variational execution is up to five times faster than the baseline approach while requiring up to 74% less memory. Beyond that, we also show that we can generate variational traces for huge configuration spaces for which it is simply impossible to generate them using a brute force approach (see Checkstyle).

To generate variational traces, we extended VaxeJ, a variational interpreter for Java [56], [57]. VaxeJ is a metacircular interpreter that executes Java bytecode instructions conditional and that stores all data using choice values. We adapted the execution of bytecode instructions, such as fields and local variable instructions to record their changes on data, if-statements to record whether they split

the execution, method invocations to record the stack trace, and exceptions to report faults.

While variational execution has been introduced before [56], [60], our new contribution is to realize its potential for generating variational traces without the disadvantages (e.g., memory overhead) of aligning single traces, by observing internals of the variational execution engine. Thus, with variational traces, we are usually able to align an exponential number of traces, which is practically impossible without it.

Generation with Symbolic Execution. Variational execution shares similar ideas with *symbolic execution* [17], [39]. Indeed, with symbolic execution it is possible to explore the executions for many inputs. In contrast, to symbolic execution, variational execution always processes alternative but concrete values, not symbolic ones. Thus, variational execution does not have the problems of symbolic execution, such as decidable loop bounds. As variational execution requires concrete inputs, our approach also requires a test case, which is given by our scenario as we want to compare executions for a given test. Symbolic execution can hardly be used to generate variational traces as it typically does not share and align executions beyond common prefixes [14], [27]. MultiSE [73], which incorporates ideas from variational execution, introduces summary values to represent value differences, which enable MultiSE to increase the sharing abilities of symbolic execution. MultiSE could potentially be used to generate variational traces as the execution may be similar to variational execution and the data-flow differences can be observed in the summary values. However, due to the overhead of the symbolic execution engine, it remains unclear whether MultiSE would achieve the same scalability and performance as variational execution.

3.4 Varviz

We argue that variational traces aid programmers to debug and understand programs by providing information about the program execution and interactions among options. To make variational traces accessible, we implemented an Eclipse plug-in called *Varviz* (from *variation* and *visualization*). We released Varviz as open-source (<https://meinicke.github.io/varviz/>). The plug-in already comes with the utilized *VarexJ* to generate variational traces.

In Figure 2, we show a screenshot of Varviz. The variational trace can be generated using default run mechanisms of Eclipse, which will automatically call *VarexJ*. After running the program, the variational trace is shown in the Varviz view (shown on the right).

Navigation is one of the most time-consuming tasks during debugging [42]. Therefore, we designed Varviz to also be used as a navigation tool. By double-clicking on elements in the trace, the tool automatically displays the file and the line of the element. As shown in the screenshot, the return instruction that throws the exception is highlighted after double-clicking the exception statement in the trace.

Focus on selected interactions. In practice, many interaction faults occur among few options [1], [13], [19], [22], [46], [54], [57], [61]. To understand a certain interaction among a specific set of options, we show a, usually smaller, relevant parts of the variational trace. To this end, we provide a new *projection mechanism* to show only the statements that explain the ef-

fects of a given set of options. All other options are set to fixed values, false if possible (the valid selection might be restricted by a feature model [33]). By setting all other options to fixed values, Varviz will produce a *projection* for the interaction of the few options of interest. For example, if a fault is thrown under context $A \wedge \neg B$, we are interested in the interactions of these two options, but not in options C and D . To create a projection on the variational trace for A and B , we set the other options C and D to false (if possible) to hide the effects and interactions of these options. To remove unnecessary elements, we create the constraint $\neg C \wedge \neg D$ and evaluate the conditions of all elements of the variational trace. If the condition of the element under the context of the constraint is satisfiable, we keep the element in the variational trace, otherwise it can be removed. Finally, we evaluate the remaining elements, whether they represent differences among the executions for the options of the projection. Removing options from the trace can highly reduce its size and thus helps to understand the interactions (as we will show in Section 5), while preserving interactions that are relevant for the options of interest.

3.5 Limitations of Variational Traces

Variational traces inherit limitations from related automatic debugging techniques based on trace alignment [81], [90]. Similar to these techniques, we compare separate executions to explain causes of faults and interactions. The similarity of these executions determines the quality of variational traces. For example, executions (i.e., test cases) can introduce minor changes (noise) that are irrelevant to the fault of interest. Thus, executions that minimize noise are always preferable. In contrast, inputs need to trigger similar executions to reveal enough information about fault and its cause. If the executions are too different, then the variational trace cannot provide enough information. When we apply variational traces to configurable systems for the same test case, the executions of the configurations will be similar by design. However, if a fault is not caused by an interaction, but simply because certain code is executed, then we can only report the context and location of the fault but not necessarily its cause.

We use the state-of-the-art variational execution engine *VarexJ* to generate variational traces [56], [57]. However, *VarexJ* has engineering limitations inherited from the interpreter of *JavaPathfinder* [26], such as incomplete support for native methods [74], multi-threading, and performance overhead due to interpreting Java bytecode. Overall, *VarexJ* is mature enough, including complex language features such as reflection, to be able to execute several medium sized Java programs [57], but not larger industrial sized programs due to the runtime overhead and its technical limitations. Variational execution is an evolving technique and advancements in variational execution will also improve the efficiency and applicability of our approach.

4 USER STUDY

Automated debugging techniques often promise large effects for debugging tasks. Previous evaluations on approaches based on execution comparison focused on reporting the size of the explanations (number of statements) instead of showing weather and how helpful they actually are for

Program	LOC	Cov	Opt	Conf	M	N	D	Instr
GameScreen	32	32	3	8	4	12	6	230
Elevator	259	193	6	20	7	12	3	5,688
NanoXML	1000	331	1	2	18	21	4	42,138

Fig. 3: Statistics on the programs used in the user study (Cov: Covered LOC by the test case, Instr: instructions executed for the test case, M: Methods, N: Nodes, D: Decisions).

debugging [81], [90]. However, only reporting quantitative results of the approach can be misleading and the expectations may not meet the reality (e.g., the explanations may be too complex and complicated to be understandable or do not contain the necessary information to understand the fault) [63]. This is the first user study on delta debugging like approaches that we are aware of.

We designed variational traces to help users to understand variations in executions. In our evaluation, we investigate how and why variational traces help users. Specifically, we perform a user study to answer the following research questions:

RQ1: *How much do variational traces improve the performance of solving debugging tasks compared to a standard debugger?* To answering RQ1, we explore the speedup and the success rates for solving debugging tasks.

RQ2: *How do variational traces help understanding differences in executions?* With RQ2, we investigate what the information needs are during a debugging task and whether the variational trace can answer them. We want to evaluate whether the provided information (i.e, the statements shown in the variational trace) is sufficient to help understanding the interactions.

Systems and Tasks. We use three subject systems in our evaluation, namely GameScreen, Elevator and NanoXML. Statistics on the programs are shown in Figure 3. We carefully chose these systems for different reasons:

GameScreen was previously used in a study which conducted the effect of different degrees of variability on program comprehension [58]. The program is a code snippet inspired by a real variability bug in BestLap, a configurable race game. Melo et al. have shown that the small program with only three features takes on average ten minutes to debug without tool support [58]. The program contains a fault that is triggered by the interaction of two features. The task is to understand the cause of the fault and the configuration in which the fault appears. The program is too trivial and cannot give any insights for our study as the program can be understood in few minutes using a standard debugger. Thus, we used GameScreen only as warm-up task to make the users familiar with the type of tasks they should perform.

Elevator is a simulation of a configurable elevator system [64]. The program is designed to trigger interactions among its options. Even though the program has only few lines of code, it is hard to understand the impact of its features due to the interactions. The program comes with several specifications in form of runtime assertions that are violated for certain configurations. We selected a specification that states that the elevator should continue in its current direction if there are still calls in this direction.

This specification is violated if a feature for *executive floors* is on, which can force the elevator to change its direction. In the tasks for Elevator, the participants should figure out in which configurations the fault appears and how the fault is caused. Although fixing a fault is part of debugging, fixing itself is not part of the task as this would have required to change the program’s specification. Instead it was sufficient to explain how the fault is caused.

NanoXML itself is not a configurable system. The program was used in a prior study to evaluate whether the automatic debugging technique Tarantula can help programmers with debugging [63]. Tarantula showed only minor improvements for debugging NanoXML compared to a standard debugger. We evaluate NanoXML on the same bug as in the original study, in line with the original study [63]. We provide two slightly different files as input for parsing. One of the files cannot be parsed correctly, causing an exception. The other one is a similar file that can be parsed. Both files are parsed simultaneously using variational execution. In addition to the tasks of the previous programs, the participants were also asked to fix the bug similar to the prior study [63]. With NanoXML we show that variational traces are helpful for a standard debugging tasks to understand variations beyond configurable systems. Thus, with the NanoXML experiment we can show the usefulness of comparative- and delta-debugging approaches which have not been evaluated in user studies before [81], [90].

Pilot Study. We performed a pilot study to estimate the required power (i.e., number of participants) of our study and to tune the task and descriptions. We asked several graduate students to use Varviz and the Eclipse debugger on our tasks. We found and revised several issues of the usability of Varviz. We also measured the time and estimated that the effect size was big enough to show significant effects with few participants in the actual experiment. For Elevator we had an estimation of 40 minutes when using the Eclipse debugger, compared to an estimation of 12 minutes when using Varviz. For, NanoXML we have an estimation of 22 minutes from a previous study when using a debugger, which we use as estimation when using a standard debugger (our results were slightly higher) [63]. We have an estimation based on our pilot study of only 8 minutes when using Varviz.

Study Design. We designed our experiment as between-subject study to compare performances between participants using a standard debugger (baseline) or Varviz (treatment). We did not mix the participants between using the standard debugger and Varviz (i.e., within-subject design). Each participants solved all three tasks with the same tool to reduce training time required for Varviz and to avoid carryover effects, such as learning effects and demand effects [16]. Learning effects from the first tasks might be applied to the second which influences the performances when using different tools. Also, the motivation of using a new tool can influence the performance of the participants. This effect is amplified if they are using both tools in a within-subject design [16]. A between-subject design has less statistical power compared to a within-subject design (i.e., we may need more participants to show significant effects), however, as we expect the effect size to be very large, as suggested by the literature on user studies [16], a between-subject design is more appropriate

as it avoids confounding factors of within-subject designs.

We did not design two comparable tasks, but intentionally two very different ones for external validity. A within subject design typically requires multiple similar tasks, which is a benefit using a between-subject design. While the participants worked on the tasks, with their consent, we recorded the screen and asked them to verbalize their thoughts using think-aloud protocols [8]. These recordings help us to track the participants' information needs and debugging strategies.

Other approaches, such as delta debugging [90] and comparative causality [81] may give similar textual explanations of the faults. However, we cannot compare our approach with delta debugging [90] and comparative causality [81] as the tools are not available (we contacted the authors) and as they are designed to explain differences among only two executions.

Methods. To answer RQ 1, we compare the time and success rates of the participants for solving the tasks. To answer RQ 2, we record the audio and the screen of the participants. We analyze the recordings based on qualitative content analysis using open coding [69], [72]. We watch the videos with the goal to find common tasks that the participants perform during debugging. We use these commonalities to create a coding frame that allows us to understand how the participants perform when using Varviz or the standard debugger.

Participants. As we plan to perform think-aloud protocols, analyzing the data (i.e., screen recordings and audio) requires high effort. We thus aim to avoid an unnecessary high number of participants. According to Nielsen [62], for think aloud protocols five participants are sufficient to gain most insights. Adding more participants does not give more essential information. Thus, to answer RQ 2, we require at least five participants per group, so ten participants in total.

To answer RQ 1, we calculate the minimum number of number of participants required using power statistics based on the pre-study results. We use the Rosner's equation to calculate the required sample size n for each group in our study [67]:

$$n = \frac{(\sigma_1^2 + \sigma_2^2)(z_{1-\alpha/2} + z_{1-\beta})^2}{\Delta^2} \quad (1)$$

We use a conservative α value of 0.01 which is a probability of Type I error of 1% (i.e. the probability to find an effect even though there is none, usually 0.05). We use a conservative β value of 0.05 which is a probability of Type II error of 5% (i.e. the probability not finding an effect even though there is one, usually 0.2). The statistical power is $1 - \beta$ and thus 95%. We use a high value for the estimated standard deviations σ_1 and σ_2 of 5 minutes as we only used few participants in our pilot study. Using Equation 1 and our estimations of expected performances (10 versus 40 and 8 versus 22 minutes for Elevator and NanoXML respectively), we can calculate the number of participants needed for our experiment as one participant ($n_{elevator} = 0.989 \approx 1$) and five participants ($n_{NanoXML} = 4.544 \approx 5$) per group for Elevator and NanoXML respectively. As the required group size for NanoXML is larger than for Elevator, we use a group size of five for our experiment. Thus, based on our pre-study results, we need ten participants to show significant results. For both, our quantitative and our qualitative analysis ten participants are sufficient.

We recruited ten participants at Carnegie Mellon University: eight undergraduate students, one graduate student and one post doc. The participants were recruited using posters and mailing lists. We excluded participants without prior knowledge of Java. All participants received a 25-dollar gift card after finishing the experiment. The participants were assigned randomly to two groups of five. One group used the Eclipse debugger, the other group used Varviz. The graduate student used Varviz, while the post doc used the Eclipse debugger.

Before conducting the experiment, we asked the participants for their programming experience and experience with Java. The groups were roughly similar: The median experience for programming is 3.5 years for both groups (average is 3.7 years for the debugger group and 5.3 years for the Varviz group; note that the large difference in average experience is caused by as a single outlier (Varviz 1) who reported 12 years of non-professional programming experience.) The median Java experience is 2 years for both groups (average is 2.6 years for the debugger group and 2 for the Varviz group). None of the participants knew variational execution or any of the subject programs. All participants in the debugger group have used Eclipse and the debugger before. If a participant did not remember how to get to a certain view, such as the call hierarchy, or were unsure about certain functionalities of Eclipse, we provided this information during the experiment.

Execution. Both groups were given an Eclipse containing the three programs they had to debug including a failing configuration for Elevator and the two XML files for NanoXML. The participants using Varviz were introduced to the functionalities of the tool. We used a simple foo-bar example to explain the functionalities of Varviz. All participant started debugging the programs in the same order, first GameScreen, second Elevator and third NanoXML. The participants performed the tasks for Elevator and NanoXML until they solved them correctly, until they gave up, until they reached a time limit of 30 minutes per task, as we planned the experiment to take roughly one hour.

As we performed a think-aloud protocol, we conducted the experiment with each participant in isolation. To record the audio and screen, we used an inhouse recording tool from our university. For the first participant using the debugger (Debugger 1), we used a different open source recording tool. We lost the screen recording due to a fault in the recording tool. However, we kept the results and the audio recording from this participant in our study as he had a good performance compared to the others using the debugger.

Quantitative Analysis. To answer RQ1, we compare the time and success rates of the participants for solving the tasks. In Figure 4, we plot the performances of the participants. The signs \checkmark and \times indicate correct and aborted solutions respectively. When using the Eclipse debugger only two out of five participants solved the Elevator task correctly and four out of five for NanoXML. In contrast, all participants using Varviz solved the tasks for both programs.

For Elevator the participants using Varviz took on average 12 minutes while the participants using the debugger took on average 28 minutes. The best performance using the debugger took 20 minutes, more than five minutes longer

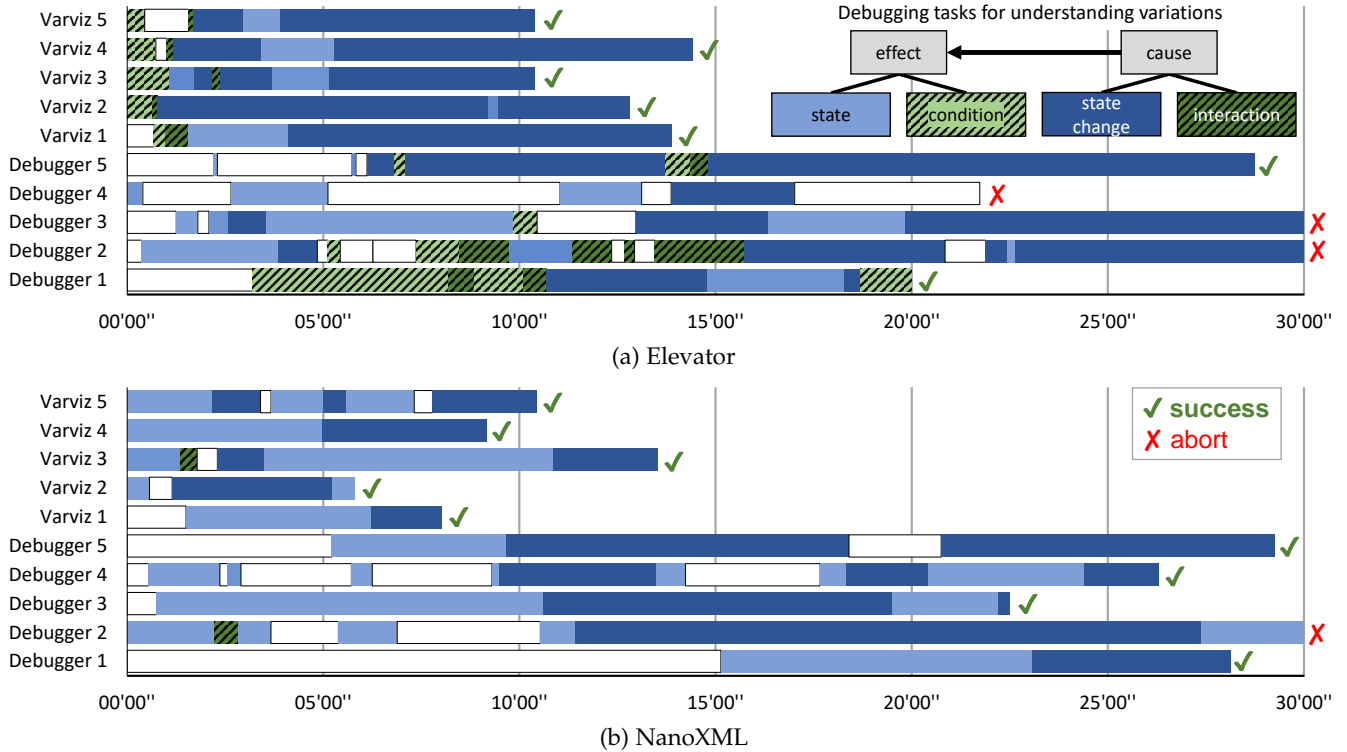


Fig. 4: Time spend on debugging tasks. White boxes denote unrelated tasks.

than the worst participant using Varviz. For NanoXML we see similar results. The Varviz group took on average 9 minutes. In contrast, the debugger group took on average 27 minutes (best was 22 minutes). Our results are in line with prior research which reported an average time of 22:30 minutes using only a debugger [63]. To calculate the *effect size* of using Varviz compared to a standard Debugger, we use the corrected equation of *Hedges' g* which corrects for the upwards bias on small sample sizes:

$$g = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2 + \sigma_2^2}{2}}} \times \frac{N - 3}{N - 2.25} \times \sqrt{\frac{N - 2}{N}} \quad (2)$$

A *g* value of 1 indicates that the groups differ by 1 standard deviation. In general, a *g* value larger than 0.8 indicates a large effect size. We calculated a *g* values of 3.6 for Elevator and 4.9 for NanoXML. Thus, *g* values indicate a huge effect of using Varviz compared to a standard Debugger for understanding differences among executions. The differences are statistically significant (Mann-Whitney U test: $p < 0.01$).

RQ1: How much do variational traces improve the performance of solving debugging tasks compared to a standard debugger?

The participants using Varviz are on average 55.4%, respectively 65.5% faster than participants using the Eclipse debugger for Elevator respectively NanoXML. The success rates when using the Eclipse debugger are 40% for Elevator and 80% for NanoXML. When using Varviz, the success rates are 100% for both programs.

Qualitative Analysis. To answer RQ 2, we analyzed the audio and the screen recordings to understand how variational traces help understanding differences in executions. Our coding frame [69], [72] is shown in the right upper corner of Figure 4a. Overall, all participants try to understand the relationship between *effect* (e.g., the fault) and its *cause* [49]. The tasks for understanding effect and cause can be refined further.

To understand the effect, it is necessary to understand program state and the condition. The program state is necessary to know values of variables and which method calls are important. When analyzing variations, it is also important to know under which condition (i.e., selection of options) the fault appears.

After understanding the effect, it is possible to investigate its cause. It is necessary to understand the state changes and method calls that lead to the state of the effect. As the effect only happens under certain conditions it is also necessary to investigate how specific selections of options cause the effect.

In Figure 4, we used our coding frame to illustrate on which tasks the participants were working on. Additionally, to the tasks in our coding frame, the participants also spend time with unrelated tasks, such as reading, scrolling, or investigating unrelated code [42]. Unrelated tasks are shown with white boxes. We plot the tasks on a horizontal time axis. The group using Varviz performed the tasks much better than the group using the debugger. In the following, we investigate the reasons why the tasks are so difficult when using only a debugger. We explore which information are required to solve the tasks and how Varviz helps to gather them.

We can see in Figure 4 that the debugger group spend much more time on unrelated parts of the program and on tasks that do not lead to solving the problem. The main reason is that the participants read unrelated code of the

programs. Another reason is that the participants give up on their current goal and try to get information from other places. When using the debugger, it is up to the programmer to find the places where to find information about the program. Thus, the participants were lost in the source code they did not know, which leads to confusion and reading of code. In contrast, when using Varviz, the participants had a guide that helps them to find the few locations in the code that are of interest.

When we analyze the performances for Elevator more closely (see Figure 4a), we see that the Varviz group took only little time to understand the variations. All the participants almost instantly figured out the condition of the fault (which is trivial using the variational trace as it is indicated by the context of the exception statement). Also, figuring out where the option affects the behavior of the program is simple using Varviz as there are only few place (few decisions) where the option affects the data flow. In contrast, when using the debugger, the first task is to figure out the condition of the fault. Without specialized tool support, this requires switching the options in the configuration and to re-execute the program. The time spent on this task depends on how many options the program has, how many options interact, and finally on luck or intuition as for participant *Debugger 5*. A simple tool that reports the condition of the faults (e.g., brute force) would help the debugger users and would improve their performances. However, such a tool alone is not sufficient as it solves only a small part of the problem.

After finding the condition of the fault, the participants still need to answer how the option triggers the effect. By searching where the option is used, the place can be found, however also unrelated usages and only direct usages are found. Thus, participants *Debugger 3* and *Debugger 4* did not even identify this part of the program.

Finally, we can see that the Varviz group spend little time to understand the state at the exception. This means that they can spend more time for understanding how the interaction triggers the effect and how this causes the fault. The debugger group took overall more time to find the values of the variables and their values at the exception (except of participant *Debugger 5* who performed well on this task). The tasks to identify the exception state and the state changes are particularly hard as the program calls the scheduling method for the elevator in a loop. This makes setting breakpoints hard as they are triggered multiple times before the actual state of interest.

In the performances for NanoXML (see Figure 4b), we can see that both groups struggle for identifying the state of the exception and answering why this state causes the fault. This is because of a relatively difficult if-statement shown in the listing below. The variational trace provides the values used in this if-statement. However, the participants still need to understand the meaning of it.

```

1 if (!str.equals(prefix==null?name:prefix+name))
2   XMLUtil.errorWrongClosingTag(this.reader, name, str);

```

After the participants using Varviz understood the meaning of the if-statement, they only spend little time to find the place of the cause as it is pointed out by a decision in the trace. In contrast, the debugger group again struggled to identify the cause. One reason is that the parsing is implemented

using recursion, which again means that the breakpoints are triggered multiple times at the wrong state. This again shows that control-flow (i.e., the decision of the cause) and data-flow information (i.e., the values that differ among the two executions) are essential to understand the differences among executions. Both information are contained in the variational trace.

RQ2: *How do variational traces help understanding differences in executions?*

Understanding where and how configuration options influence the execution is hard using only a standard debugger as only one configuration can be executed at a time. Thus, participants using the Eclipse debugger have difficulties figuring out the **condition** of the fault, gathering information about differences in the **program states**, and to find the **cause** of the fault. Variational traces help with these task by providing essential information about the **fault condition**, the **fault state**, as well as **data and control-flow differences** that lead to the fault. This information helps users to focus on important parts of the execution which additionally prevents from wasting time on unrelated activities.

Threats to Validity. We designed our user study as between-subject study with only ten participants. As discussed, we decided to use a between-subject design to avoid confounding factors and to reduce training time at the cost of the reduced power of the design [16]. We carefully calculated the required number of users upfront to minimize the effort for analyzing the think-aloud protocols and the screen recordings [67]. The measured performances in the experiment approximately match our expected performances from our pilot study. Due to the large effect sizes a larger number of participants is unnecessary. The fact that our results are statistically significant confirms that the chance of a random error due to small sample size is small.

To reduce selection bias, we recruit participants in public channels and randomly assign them into two groups. The average programming experience varies by almost two years between the groups, which is however caused by a single outlier (Varviz 1). Our results are robust to removing this outlier (i.e., our results remain statistically significant without this participant) [48]. To avoid effects due to differences in programming experiences, we designed the tasks in a way that basic debugging experience is sufficient. Since we conduct the user study in the Eclipse environment, experience with the Eclipse toolchain is likely to affect performance of participants. We performed a warm-up task familiarize the participants with the type of tasks and the programming environment. We argue that most common usages of Eclipse are straightforward to most developers, given that Eclipse is a standard and classic development tool for Java. In addition, we made it clear to the participants before study that we could provide immediate support for questions on Eclipse usage. However, participants rarely asked for help regarding Eclipse. Thus, we argue that Eclipse experience likely has only minor impact to the performance results. To minimize confounding factors, we implement Varviz in a way that

is completely orthogonal with existing features of Eclipse. Moreover, our introduction to Varviz only covers the plugin itself, not including any other functionalities of Eclipse. The think-aloud protocol influences the time performance of the participants; however, it influences both groups equally and because the expected effect size is big we do not expect any systematic influences on the overall results.

We used two diverse systems for the debugging tasks. We showed that variational traces are useful to understand faults in systems with multiple options as well to compare two executions. However, readers should be careful when generalizing our results to other systems and tasks.

5 SCALABILITY EVALUATION

Variational traces are concise representations of differences among executions. To further reduce the size, we allow focusing on small sets of options discussed in Section 3.4. However, we do not yet apply any kind of impact analyses to reduce the size even further, as this is out of scope of this paper and as the sizes are already small enough especially for the programs used in Section 4. Variational traces are useful beyond debugging, as for example in our work on detect behavioral feature interactions with feature interaction graphs which can deal with large variational traces [75]. In this section, we evaluate the size of variational traces when aligning the executions of exponentially large configuration spaces. Specifically, we answer the following research question:

RQ3: *How does the generation of variational traces using variational executions scale compared to a base line approach?* The exponential growth of configuration spaces with the number of options is challenging for both, execution and alignment of many configurations. By answering RQ 3, we investigate the scalability of using variational execution to generate variational traces with regard to runtime and memory consumption.

RQ4: *How large do variational traces get?* Information on data and control flow differences are beyond debugging (cf. Section 6). With RQ4, we investigate how complex variational traces get when applied to programs with different numbers of options and different sizes. We also want to find out how effective the filters for options are for reducing the size (see Section 3.4).

Experimental Setup. In our evaluation, we reuse six configurable systems from our previous study on feature interactions and essential configuration complexity [57]. The systems are shown in Figure 3. The systems are from different domains and show different interaction properties [57]. We execute each system for a corresponding standard scenario. Each system comes with a set of options that can be enabled and disabled, which results in large numbers of configurations for which we execute the program (see Figure 5).

The experiments are performed on a Windows computer with 8 GB ram and an Intel i5 processor with 4 cores. To answer RQ 3 we collect the execution time and the memory consumption, we ran variational execution and the base line approach ten times and report the median value to avoid measurement errors.

The size of the variational trace depends on several factors. First, the length of the execution (see number of instructions in Figure 5). It also depends on how the program implements variability and how the options interact. The fewer options interact in a program, the shorter the trace will be (statistics on the interactions in the systems are discussed by Meinicke [57]). In our evaluation, we collect metrics on how large variational traces can get for different implementations and executions.

Statements and decisions indicate the complexity of the variational trace. However, understanding a fault or an interaction usually only requires few options as the interactions degrees are usually low [1], [13], [19], [22], [46], [54], [57], [61]. We use context filters as discussed in Section 3.4 to filter the variational trace for all combinations of three options (we filter the trace of CheckStyle only for *one* option due to the large number of options).

Results. In Figure 5, we report the times and memory consumptions required to generate the variational trace. As shown, the time to generate the variational trace is always lower than with the base line approach, especially for larger configuration spaces. For Checkstyle with 141 options, our approach takes 209.8 seconds, which is lower than what the base line approach takes when aligning the configurations for only five options.

As expected, the memory consumption of the base line approach becomes problematic, especially when aligning longer execution traces. Again, the base line approach requires 3 GB for aligning only the traces for five options while our approach takes 2 GB when aligning the traces for *all* configurations. Note that the reported memory consumption of our approach also contains the memory overhead of VaxexJ itself.

RQ3: *How does the generation of variational traces using variational executions scale compared to a base line approach?*

Generating variational traces using variational executions scales to exponentially large configuration spaces with regard to execution time and memory consumption. In contrast, the base line approach is only able to generate variational traces for smaller configuration spaces while taking more time and memory than our approach based on variational execution.

In Figure 5, we report the sizes of complete variational traces and the mean sizes after applying filters for three options. Even though the interaction experiments (Elevator, E-Mail and Mine Pump) are designed to cause many interactions among options, we see that the sizes of their variational traces are small. For GPL the complete variational trace becomes large as the executions contain long iterations which cause trivial but repetitive interactions on data (e.g., optionally initializing the weight for all vertexes in a graph). The same happens for QuEval which also has trivial but repetitive executions. For CheckStyle, which has by far the longest executions, we see that the size of the complete variational trace is huge. The size of the variational trace is, however, small in relation to the configuration space and the number of executed instructions. Most of these statements in the variational trace for CheckStyle are again repetitive.

Program	SLOC	Opt.	Conf.	Instructions	Time	TimeBL	Memory	MemoryBL	D_{all}	D_3	S_{all}	S_3
CheckStyle	14,950	141	$> 2^{135}$	194,725,919	209.8s	*364.8.8s	1984MB	*3269MB	5,989	165.50	290,477	2022.7
QuEval [70]	3,109	23	940	6,460,571	10.6s	50.3s	379MB	1498MB	699	26.9	4,152	110.0
GPL [52]	662	15	146	17,457,437	15.6s	18.9s	408MB	975MB	530	6.9	5,565	301.0
Elevator [64]	730	6	20	23,559	0.1s	0.2s	41MB	29MB	36	17.7	96	51.0
E-Mail [25]	644	9	40	25,846	0.2s	0.4s	28MB	49MB	53	9.7	129	17.5
Mine Pump [45]	296	6	64	21.615	0.1s	0.4s	37MB	49MB	10	7.6	14	10.9

Fig. 5: Statistics on programs used in quantitative evaluation. D_{all} and S_{all} state the number of decisions respectively statements of the full variational trace. D_3 and S_3 state the mean number of decisions respectively statements after filtering for three options. In TimeBL and MemoryBL we show the time and memory consumptions for the baseline implementation (*for CheckStyle we only executed the configurations for five options with the baseline approach).

After applying the filter for one option, we can see that the size can be reduced a lot. In general, the sizes of the traces become small after filtering them. Even the full traces can be useful as most of the shown statements and decisions are repetitive due to iterations.

RQ4: *How large do variational traces get?*

Variational traces are often small, but can become large, especially due to long iterations that repetitively create the same interactions on data. The size can be drastically reduced due to the filtering mechanisms, such as filtering for a small set of options.

Threats to Validity. To mitigate threats to external validity, we analyze different programs that show different kinds of interactions. We argue that variational traces are scalable (i.e., we can align the execution for an exponential configuration space while the number of nodes does not grow exponentially with the number of configurations) to most programs in the wild, because recent studies have shown that although there could be many options in a program, most options interact locally and thus interaction degrees are usually low [1], [13], [19], [22], [46], [54], [57], [61].

6 APPLICATIONS BEYOND DEBUGGING

Understanding differences among executions can be useful beyond debugging. In this section, we discuss three further potential applications of variational traces.

Detecting behavioral feature interactions. Detecting faults requires certain types of specifications (e.g., in form of test cases). However, certain types of feature interaction faults are hard to detect if they do not trigger an exception but only have behavioral impacts. For example, two features may interact in such a way that one feature suppresses the effect of another. In previous work [75], we presented *feature interactions graphs* that help to identify which feature interact with each other and whether their interactions are suspicions (e.g., suppression). These feature interactions graphs are based on an analysis of variational traces.

Understanding impact of load-time options. It is hard to identify code that is affected by load-time options, especially due to implicit data flow. Previous work used static taint analysis to detect code that depends on the selection of certain options [51]. However, it is hard to trace back why

the code is affected as such information is lost in the analysis. With variational traces we can help to understand why certain parts of the code depends on the selection of options as we show causes of differences in the control flow.

Understanding information flow. Previous work compared executions to detect information flow, either using a similar technique to variational execution [6], [7], [71], [88] or using multi execution [20], [38], [43], [47]. Aligning executions allows to detecting potential leaks of secret data. However, to understand why the information is leaked requires understanding why the executions differ. Thus, again, variational traces can help understanding the causes of information flow by tracing data differences.

7 RELATED WORK

We already discussed closely related work in the domains of automatic debugging [2], [4], [23], [31], [31], [63], [79], [81], [87], [90], [91] and feature interactions [19], [35], [37], [51], [54], [57], [60], [61], [76], [83] in Section 2. Our work combines ideas from both fields and builds specifically on the idea of sharing and coordinating multiple executions with variational execution [56], [57], [60] and thus sidesteps the challenges of trace alignment and full trace recording [12], [32] as discussed.

Omniscient or *back-in-time debugging* allows exploring and debug a single execution [11], [40], [41], [50], [66]. In contrast to standard debuggers, back-in-time debuggers allow exploring information of previous parts of the execution. To do so, they record full traces resulting in severe scalability challenges as they cannot predict which information is of interest. In contrast, we provide a dynamic analysis that can decide on the fly which information is potentially relevant to explain the fault.

Symbolic execution allows to explore all executions of a program for different inputs. As discussed in Section 3.3, symbolic execution usually does not share the executions after they are separated unless they incorporate ideas from variational execution [73]. The *interactive verification debugger* is a tool to understand the symbolic execution of a program [27], [28]. The tool visualizes the execution in a tree structure similar to Varviz. However, as the symbolic execution never joins the tree structure gets large even for small programs. In contrast, our variational trace provides a concise representation of many executions.

Static and dynamic program slicing are techniques to reduce a program to only the statements relevant for understanding

the state at a given program point [44], [85]. However, program slices are often large and cannot explain differences among multiple executions, especially execution omission bugs [94]. Differential slicing and dual slicing [30], [38], [81] compare the execution of two executions and reduces the comparison using program slicing. In contrast to slicing approaches, we aim to explain the differences among many executions and only focus on the causes and differences in the state to keep the explanations concise.

Multi execution are approaches that synchronize (typically two) concrete executions. These approaches enable analyses for information flow [20], [38], [43], [47], configuration faults [78] and inconsistent updates [29], [53], [82]. In contrast, our approach can compare a potentially exponential number of executions and helps to understand how the differences affect the program behavior.

8 CONCLUSION

In this work, we propose variational traces to explain the runtime behaviors of inputs and interactions among them. Using variational execution and specialized filters, we scale the generation of variational traces to the potentially exponential space of possible inputs. To visualize variational traces, we provide an interactive Eclipse plugin called Varviz, which enables programmers to use variational traces for debugging interaction faults. In our user study, we show that users of Varviz outperform the users of the Eclipse debugger significantly in terms of understanding and time spent on debugging tasks. Users of Varviz can focus on relevant parts of the programs quickly, without being distracted by irrelevant data and control flow decisions. When compared with users who use standard Eclipse debugger, Varviz users can finish all the debugging and understanding tasks, using less than half of the time. We further evaluate the size of variational traces on six highly configurable systems. In general, the size of variational traces can get large, but our filters are effective in reducing the traces to a relatively small number of statements. Overall, our evaluation of effectiveness and scalability demonstrates that variational traces are useful in practice to understand differences among executions.

REFERENCES

- [1] I. Abal, C. Brabrand, and A. Wasowski. 42 Variability Bugs in the Linux Kernel: A Qualitative Analysis. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 421–432. ACM, 2014.
- [2] R. Abreu, P. Zoetewij, and A. J. C. Van Gemund. An evaluation of similarity coefficients for software fault localization. In *Proceedings of the Pacific Rim International Symposium on Dependable Computing (PRDC)*, pages 39–46, 2006.
- [3] S. Apel, S. Kolesnikov, N. Siegmund, C. Kästner, and B. Garvin. Exploring Feature Interactions in the Wild: The New Feature-Interaction Challenge. In *Proceedings of the International SPLC Workshop Feature-Oriented Software Development (FOSD)*, pages 1–8. ACM, 2013.
- [4] T. Apiwattanapong, A. Orso, and M. J. Harrold. JDiff: A Differencing Technique and Tool for Object-Oriented Programs. *Proceedings of the International Conference on Automated Software Engineering (ASE)*, 14(1):3–36, 2007.
- [5] T. H. Austin and C. Flanagan. Multiple Facets for Dynamic Information Flow. *ACM Sigplan Notices*, 47(1):165–178, 2012.
- [6] T. H. Austin and C. Flanagan. Multiple Facets for Dynamic Information Flow. In *Proceedings of the Symposium on Principles of Programming Languages (POPL)*, pages 165–178. ACM, 2012.
- [7] T. H. Austin, J. Yang, C. Flanagan, and A. Solar-Lezama. Faceted Execution of Policy-Agnostic Programs. In *Proceedings of the ACM SIGPLAN Workshop on Programming Languages and Analysis for Security (PLAS)*, pages 15–26. ACM, 2013.
- [8] H. Beyer and K. Holtzblatt. *Contextual Design: Defining Customer-Centered Systems*. Elsevier, 1997.
- [9] M. Böhme, B. C. d. S. Oliveira, and A. Roychoudhury. Regression Tests to Expose Change Interaction Errors. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 334–344. ACM, 2013.
- [10] J. R. Burch, E. M. Clarke, K. L. McMillan, D. L. Dill, and L.-J. Hwang. Symbolic Model Checking: 10²⁰ States and Beyond. In *Symposium on Logic in Computer Science (LICS)*, pages 428–439. IEEE, 1990.
- [11] B. Burg, R. Bailey, A. J. Ko, and M. D. Ernst. Interactive Record/Replay for Web Application Debugging. In *Proceedings of ACM symposium on User interface software and technology (SIGHCI)*, pages 473–484. ACM, 2013.
- [12] M. Burtcher, I. Ganusov, S. J. Jackson, J. Ke, P. Ratanaworabhan, and N. B. Sam. The vpc trace-compression algorithms. *IEEE Journal of Transactions on Computers (TC)*, 54(11):1329–1344, 2005.
- [13] I. Cabral, M. B. Cohen, and G. Rothermel. Improving the Testing and Testability of Software Product Lines. In *Proceedings of the International Software Product Line Conference (SPLC)*, pages 241–255. Springer, 2010.
- [14] C. Cadar, P. Godefroid, S. Khurshid, C. S. Păsăreanu, K. Sen, N. Tillmann, and W. Visser. Symbolic Execution for Software Testing in Practice: Preliminary Assessment. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 1066–1071. ACM, 2011.
- [15] M. Calder, M. Kolberg, E. H. Magill, and S. Reiff-Marganec. Feature Interaction: A Critical Review and Considered Forecast. *Computer Networks*, 41(1):115–141, 2003.
- [16] G. Charness, U. Gneezy, and M. A. Kuhn. Experimental Methods: Between-Subject and Within-Subject Design. *Journal of Economic Behavior & Organization*, 81(1):1–8, 2012.
- [17] L. A. Clarke. A Program Testing System. In *ACM*, pages 488–491. ACM, 1976.
- [18] A. Classen, P. Heymans, P.-Y. Schobbens, and A. Legay. Symbolic Model Checking of Software Product Lines. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 321–330. ACM, 2011.
- [19] M. B. Cohen, M. B. Dwyer, and J. Shi. Interaction Testing of Highly-Configurable Systems in the Presence of Constraints. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA)*, pages 129–139. ACM, 2007.
- [20] D. Devriese and F. Piessens. Noninterference Through Secure Multi-Execution. In *Symposium on Security and Privacy (SP)*, pages 109–124. IEEE Computer Science, 2010.
- [21] M. Erwig and E. Walkingshaw. The Choice Calculus: A Representation for Software Variation. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 21(1):6:1–6:27, 2011.
- [22] B. J. Garvin and M. B. Cohen. Feature Interaction Faults Revisited: An Exploratory Study. In *Proceedings of the International Symposium on Software Reliability Engineering (ISSRE)*, pages 90–99. IEEE Computer Science, 2011.
- [23] A. Groce, S. Chaki, D. Kroening, and O. Strichman. Error Explanation with Distance Metrics. *International Journal on Software Tools for Technology Transfer (STTT)*, 8(3):229–247, 2006.
- [24] Z. Gu, E. T. Barr, D. J. Hamilton, and Z. Su. Has the Bug Really Been Fixed? In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 55–64. ACM, 2010.
- [25] R. J. Hall. Fundamental Nonmodularity in Electronic Mail. *Proceedings of the International Conference on Automated Software Engineering (ASE)*, 12(1):41–79, 2005.
- [26] K. Havelund and T. Pressburger. Model Checking Java Programs Using Java PathFinder. *International Journal on Software Tools for Technology Transfer (STTT)*, 2(4):366–381, 2000.
- [27] M. Hentschel, R. Hähnle, and R. Bubel. An Empirical Evaluation of Two User Interfaces of an Interactive Program Verifier. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 403–413, 2016.
- [28] M. Hentschel, R. Hähnle, and R. Bubel. The Interactive Verification Debugger: Effective Understanding of Interactive Proof Attempts. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 846–851, 2016.
- [29] P. Hosek and C. Cadar. Safe Software Updates via Multi-Version Execution. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 612–621. IEEE Press, 2013.

- [30] N. M. Johnson, J. Caballero, K. Z. Chen, S. McCamant, P. Poosankam, D. Reynaud, and D. Song. Differential Slicing: Identifying Causal Execution Differences for Security Applications. In *Symposium on Security and Privacy (SP)*, pages 347–362. IEEE, 2011.
- [31] J. A. Jones, M. J. Harrold, and J. Stasko. Visualization of Test Information to Assist Fault Localization. pages 467–477. ACM, 2002.
- [32] S. Kanev and R. Cohn. Portable trace compression through instruction interpretation. In *Proceedings of the International IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, pages 107–116. IEEE, 2011.
- [33] K. C. Kang, S. G. Cohen, J. A. Hess, W. E. Novak, and A. S. Peterson. Feature-Oriented Domain Analysis (FODA) Feasibility Study. Technical Report CMU/SEI-90-TR-21, Software Engineering Institute, 1990.
- [34] C. H. P. Kim, D. Batory, and S. Khurshid. Eliminating Products to Test in a Software Product Line. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 139–142. ACM, 2010.
- [35] C. H. P. Kim, D. Batory, and S. Khurshid. Reducing Combinatorics in Testing Product Lines. In *Proceedings of the International Conference on Aspect-Oriented Software Development (AOSD)*, pages 57–68. ACM, 2011.
- [36] C. H. P. Kim, S. Khurshid, and D. Batory. Shared Execution for Efficiently Testing Product Lines. In *Proceedings of the International Symposium on Software Reliability Engineering (ISSRE)*, pages 221–230. IEEE Computer Science, 2012.
- [37] C. H. P. Kim, D. Marinov, S. Khurshid, D. Batory, S. Souto, P. Barros, and M. d’Amorim. SPLat: Lightweight Dynamic Analysis for Reducing Combinatorics in Testing Configurable Systems. In *Proceedings of the European Software Engineering Conference/Foundations of Software Engineering (ESEC/FSE)*, pages 257–267. ACM, 2013.
- [38] D. Kim, Y. Kwon, W. N. Sumner, X. Zhang, and D. Xu. Dual Execution for On the Fly Fine Grained Execution Comparison. In *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pages 325–338, 2015.
- [39] J. C. King. Symbolic Execution and Program Testing. *Communications of the ACM*, 19(7):385–394, 1976.
- [40] A. J. Ko and B. A. Myers. Debugging Reinvented: Asking and Answering Why and Why Not Questions about Program Behavior. In *Proceedings of the International Conference on Software Engineering (ICSE)*, page 301, 2008.
- [41] A. J. Ko and B. A. Myers. Extracting and answering why and why not questions about Java program output. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 20(2):4:1–4:36, 2010.
- [42] A. J. Ko, B. A. Myers, M. J. Coblenz, and H. H. Aung. An Exploratory Study of How Eevelopers Seek, Relate, and Collect Relevant Information During Software Maintenance Tasks. *IEEE Transactions on Software Engineering (TSE)*, 32(12), 2006.
- [43] C. Kolbitsch, B. Livshits, B. Zorn, and C. Seifert. Rozzle: Decloaking Internet Malware. In *IEEE Symposium on Security and Privacy (SP)*, pages 443–457. IEEE, 2012.
- [44] B. Korel and J. Laski. Dynamic Program Slicing. *Information processing letters*, 29(3):155–163, 1988.
- [45] J. Kramer, J. Magee, M. Sloman, and A. Lister. CONIC: An Integrated Approach to Distributed Computer Control Systems. *IEE Proc. Computers and Digital Techniques*, 130(1):1–, 1983.
- [46] D. R. Kuhn, D. R. Wallace, and A. M. Gallo. Software Fault Interactions and Implications for Software Testing. *IEEE Transactions on Software Engineering (TSE)*, 30:418–421, 2004.
- [47] Y. Kwon, D. Kim, W. N. Sumner, K. Kim, B. Saltaformaggio, X. Zhang, and D. Xu. LDx: Causality Inference by Lightweight Dual Execution. In *Proceedings of of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, volume 51, pages 503–515. ACM, 2016.
- [48] H. Larsson, E. Lindqvist, and R. Torkar. Outliers and Replication in Software Engineering. In *Proceedings of the Asia-Pacific Software Engineering Conference (APSEC)*, volume 1, pages 207–214. IEEE, 2014.
- [49] T. D. LaToza, D. Garlan, J. D. Herbsleb, and B. A. Myers. Program Comprehension as Fact Finding. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 361–370. ACM, 2007.
- [50] B. Lewis. Debugging Backwards in Time. *Computing Research Repository (CoRR)*, 2003.
- [51] M. Lillack, C. Kästner, and E. Bodden. Tracking Load-Time Configuration Options. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 445–456. ACM, 2014.
- [52] R. E. Lopez-Herrejon and D. Batory. A Standard Problem for Evaluating Product-Line Methodologies. In *Proceedings of the International Conference on Generative and Component-Based Software Engineering (GCSE)*, pages 10–24. Springer, 2001.
- [53] M. Maurer and D. Brumley. Tachyon: Tandem Execution for Efficient Live Patch Testing. In *USENIX Security Symposium*, pages 617–630, 2012.
- [54] F. Medeiros, C. Kästner, M. Ribeiro, R. Gheyi, and S. Apel. A Comparison of 10 Sampling Algorithms for Configurable Systems. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 664–675. ACM, 2016.
- [55] F. Medeiros, C. Kästner, M. Ribeiro, S. Nadi, and R. Gheyi. The Love/Hate Relationship with the C Preprocessor: An Interview Study. In *Proceedings of the European Conference on Object-Oriented Programming (ECOOP)*, volume 37 of *LIPICs*, pages 495–518. Schloss Dagstuhl–LZI, 2015.
- [56] J. Meinicke. VaxxJ: A Variability-Aware Interpreter for Java Applications. Master’s thesis, University of Magdeburg, 2014.
- [57] J. Meinicke, C. P. Wong, C. Kästner, T. Thüm, and G. Saake. On Essential Configuration Complexity : Measuring Interactions in Highly-Configurable Systems. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, number 2, pages 483–494, 2016.
- [58] J. Melo, C. Brabrand, and A. Wasowski. How Does the Degree of Variability Affect Bug Finding? In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 679–690. ACM, 2016.
- [59] S. B. Needleman and C. D. Wunsch. A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins. *Journal of molecular biology*, 48(3):443–453, 1970.
- [60] H. V. Nguyen, C. Kästner, and T. N. Nguyen. Exploring Variability-Aware Execution for Testing Plugin-Based Web Applications. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 907–918. ACM, 2014.
- [61] C. Nie and H. Leung. A Survey of Combinatorial Testing. *ACM Computing Surveys*, 43(2):11:1–11:29, 2011.
- [62] J. Nielsen. Estimating the Number of Subjects Needed for a Thinking Aloud Test. *International Journal of Human-Computer Studies(IJHCS)*, 41(3):385–397, 1994.
- [63] C. Parnin and A. Orso. Are Automated Debugging Techniques Actually Helping Programmers? In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA)*, pages 199–209. ACM, 2011.
- [64] M. Plath and M. Ryan. Feature Integration Using a Feature Construct. *Science of Computer Programming (SCP)*, 41(1):53–84, 2001.
- [65] G. Pothier and É. Tanter. Back to the Future: Omniscient Debugging. *IEEE software*, 26(6), 2009.
- [66] G. Pothier, É. Tanter, and J. Piquer. Scalable Omniscient Debugging. In *Proceedings of the Conference on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, volume 42, pages 535–552. ACM, 2007.
- [67] B. Rosner. *Fundamentals of Biostatistics*. Nelson Education, 2015.
- [68] J. Rubin and M. Chechik. N-Way Model Merging. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 301–311. ACM, 2013.
- [69] J. Saldaña. *The Coding Manual for Qualitative Researchers*. Sage, 2015.
- [70] M. Schäler, A. Grebhahn, R. Schröter, S. Schulze, V. Köppen, and G. Saake. QuEval: Beyond High-Dimensional Indexing à la Carte. In *Proceedings of of the International Conference on Very Large Data Bases (VLDB)*, pages 1654–1665. VLDB Endowment, 2013.
- [71] T. Schmitz, D. Rhodes, T. H. Austin, K. Knowles, and C. Flanagan. Faceted Dynamic Information Flow via Control and Data Monads. In *POST*, pages 3–23, 2016.
- [72] M. Schreier. *Qualitative Content Analysis in Practice*. SAGE Publications, 2012.
- [73] K. Sen, G. Necula, L. Gong, and W. Choi. MultiSE: Multi-Path Symbolic Execution Using Value Summaries. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 842–853. ACM, 2015.
- [74] N. Shafiei and F. v. Breugel. Automatic Handling of Native Methods in Java PathFinder. In *Proceedings of the International SPIN Symposium on Model Checking of Software (SPIN)*, pages 97–100. ACM, 2014.

- [75] L. R. Soares, J. Meinicke, S. Nadi, C. Kästner, and E. S. de Almeida. VarXplorer: Lightweight Process for Dynamic Analysis of Feature Interactions. In *Proceedings of the Workshop on Variability Modelling of Software-intensive Systems (VaMoS)*, pages 59–66. ACM, 2018.
- [76] S. Souto, M. d’Amorim, and R. Gheyi. Balancing Soundness and Efficiency for Practical Testing of Configurable Systems. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 632–642. IEEE Press, 2017.
- [77] M. Sridharan, S. J. Fink, and R. Bodík. Thin Slicing. In *Proceedings of the International Conference on Programming Language Design and Implementation (PLDI)*, pages 112–122, 2007.
- [78] Y.-Y. Su, M. Attariyan, and J. Flinn. Autobash: Improving configuration management with operating system causality analysis. In *Proc. Symposium on Operating Systems Principles*, pages 237–250. ACM, 2007.
- [79] W. N. Sumner and X. Zhang. Algorithms for automatically computing the causal paths of failures. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 355–369, 2009.
- [80] W. N. Sumner and X. Zhang. Memory Indexing: Canonicalizing Addresses Across Executions. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 217–226. ACM, 2010.
- [81] W. N. Sumner and X. Zhang. Comparative causality: Explaining the differences between executions. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 272–281. IEEE Press, 2013.
- [82] J. Tucek, W. Xiong, and Y. Zhou. Efficient Online Validation with Delta Execution. *ACM SIGARCH Computer Architecture News (SIGARCH)*, 37(1):193–204, 2009.
- [83] A. von Rhein, S. Apel, and F. Raimondi. Introducing Binary Decision Diagrams in the Explicit-State Verification of Java Code. 2011.
- [84] D. Weeratunge, X. Zhang, W. N. Sumner, and S. Jagannathan. Analyzing Concurrency Bugs Using Dual Slicing. In *ISSTA*, pages 253–264, 2010.
- [85] M. Weiser. Program Slicing. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 439–449. IEEE Computer Science, 1981.
- [86] W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa. A Survey on Software Fault Localization. *IEEE Transactions on Software Engineering (TSE)*, 42(8):707–740, 2016.
- [87] B. Xin, W. N. Sumner, and X. Zhang. Efficient Program Execution Indexing. In *Proceedings of the International Conference on Programming Language Design and Implementation (PLDI)*, pages 238–248. ACM, 2008.
- [88] J. Yang, T. Hance, T. H. Austin, A. Solar-Lezama, C. Flanagan, and S. Chong. Precise, Dynamic Information Flow for Database-Backed Applications. *arXiv preprint arXiv:1507.03513*, 2015.
- [89] A. Zeller. Yesterday, my Program Worked. Today, it Does Not. Why? In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 253–267, 1999.
- [90] A. Zeller. Isolating Cause-Effect Chains From Computer Programs. In *Proceedings of the International Symposium Foundations of Software Engineering (FSE)*, pages 1–10. ACM, 2002.
- [91] A. Zeller and R. Hildebrandt. Simplifying and Isolating Failure-Inducing Input. *IEEE Transactions on Software Engineering*, 28:183–200, 2002.
- [92] S. Zhang and M. D. Ernst. Automated Diagnosis of Software Configuration Errors. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 312–321. IEEE Computer Science, 2013.
- [93] S. Zhang and M. D. Ernst. Which Configuration Option Should I Change? In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 152–163. ACM, 2014.
- [94] X. Zhang, S. Tallam, N. Gupta, and R. Gupta. Towards Locating Execution Omission Errors. In *Proceedings of the International Conference on Programming Language Design and Implementation (PLDI)*, pages 415–424. ACM, 2007.