

One λ at a time: What do we know about presenting human-friendly output from program analysis tools?

A Scoping Review of PLDI Proceedings for HCI Researchers

Titus Barik
North Carolina State University

Chris Parnin
North Carolina State University

Emerson Murphy-Hill
North Carolina State University

Abstract

Program analysis tools perform sophisticated analysis on source code to help programmers resolve compiler errors, apply optimizations, and identify security vulnerabilities. Despite the utility of these tools, research suggests that programmers do not frequently adopt them in practice—a primary reason being that the output of these tools are difficult to understand. Towards providing a synthesis of what researchers know about the presentation of program analysis output to programmers, we conducted a scoping review of the PLDI conference proceedings. The scoping review serves as interim guidance for advancing collaborations between research disciplines. We discuss how cross-disciplinary communities, such as PLATEAU, are critical to improving the usability of program analysis tools.

ACM Reference Format:

Titus Barik, Chris Parnin, and Emerson Murphy-Hill. 2017. One λ at a time: What do we know about presenting human-friendly output from program analysis tools?. In *Proceedings of ACM SIGPLAN Conference on Programming Languages, New York, NY, USA, January 01–03, 2017 (PL’17)*, 8 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

In 1983, Brown lamented that one of the most neglected aspects of the human-machine interface was the quality of the error messages produced by the machine. Today, it appears that many of Brown’s laments still hold true with regard to program analysis tools—tools that are intended to help programmers resolve defects in their code. For example, interview and survey studies conducted at Microsoft reveal that poor error messages remain one of the top pain points when using program analysis tools [14], and other studies show similar frustration with error messages in tools [7, 30, 51]. In academia, the situation seems even more dire. As Hanenberg noted in his essay on programming languages research: “developers, which are the main audience for new language constructs, are hardly considered in the research process.” And Danas et al. note that in some cases, the output

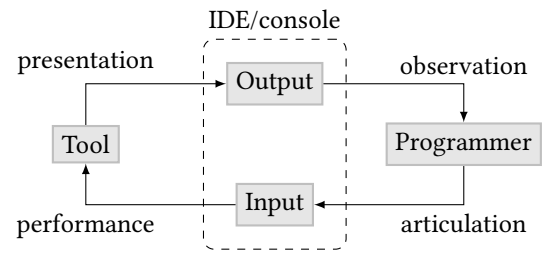


Figure 1. The interaction framework.

of program analysis tools, such as in model-finders and SAT-solvers, are generated arbitrarily and in an unprincipled way, without regard to the friendliness towards the programmer who might actually use them [17].

In prior work, we have modeled the interaction of programmers with their program analysis tools in terms of an interaction framework, conceptualized by Abowd and Beagle [1] and adapted to tools by Traver [57] (Figure 1). The interaction framework describes the different interactions between the tool and the programmer, with the tool performing some sophisticated analysis, presenting the information to the programmer through a console or IDE, and then allowing the developer to articulate their intentions back to the tool. In this paper, we are interested specifically in the *presentation* aspect of the framework, and what we know about presenting human-friendly output from program analysis tools.

Towards the longer-term goal of providing a comprehensive knowledge synthesis about program analysis output, we conducted an interim scoping review of the proceedings from Programming Language Design and Implementation (PLDI), from 1988-2017. The scoping review is intended to be accessible to human-computer interaction (HCI) researchers who want to understand how the PL community is currently applying program analysis output, in order to eventually bridge HCI research with program analysis tools. Consequently, while PLDI papers are typically written to emphasize the formal properties of their program analysis tools as their primary goal, our scoping review reframes these papers in terms of their *program analysis output* as the primary investigation.

The contributions of this scoping review are:

PL’17, January 01–03, 2017, New York, NY, USA
2017. ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

- A *quasi-gold set* of manually-identified papers from PLDI that relate to program analysis output, to bootstrap future, comprehensive literature reviews on the subject of human-friendly program analysis output.
- A knowledge synthesis of the features of program analysis output that researchers employ to present output to programmers, instantiated as a taxonomy (Section 3). Our taxonomy is agnostic to a particular mode of output, such as text or graphics.
- A bridge from HCI to PL, to foster collaboration between researchers in both communities, and to familiarize the HCI community with program analysis tools in PL.

2 Methodology

2.1 What is a Scoping Review?

In this study, we conduct a *scoping review*—a reduced form of the traditional systematic literature review [5, 32]. Scoping reviews have many of the same characteristics of traditional literature reviews: their purpose is to collect, evaluate, and present the available research evidence for a particular investigation. However, because of their reduced form, they can also be executed more rapidly than traditional literature reviews [50]. For example, reductions to scoping reviews include limiting the types of literature databases, constraining the date range under investigation, or eliding consistency measures such as inter-rater agreement. A notable weakness of scoping reviews are that they are not a final output; instead, they provide interim guidance towards what can be expected if a comprehensive literature review is conducted. Scoping reviews are particularly useful in this interim stage for soliciting guidance on conducting a more formal review, as is our intention in this paper.

2.2 Execution of SALSA Framework

We conducted our scoping review using the traditional SALSA framework: Search, Appraisal, Synthesis, and Analysis. Here, we discuss the additional constraints we adopted in using SALSA for our scoping review.

Search. We scoped our search to all papers within a single conference: Programming Language Design and Implementation (PLDI), for all years (1988-2017). As HCI researchers, we selected PLDI because it is considered to be a top-tier conference for programming languages research, because it contains a variety of program analysis tools, and because these tools tend to have formal properties of soundness and completeness that are not typically found in prototype tools within HCI. Discussions with other researchers within PLDI also revealed that researchers are interested in having their tools adopted by a broader community, but confusing program analysis output hinders usability of the tools to users outside their own research groups.

Appraisal. We manually identified papers through multiple passes. In the first pass, we skimmed titles and abstracts and included any papers which mentioned a program analysis tool and indicated output intended to be consumed by a programmer other than the authors of the tool. In this pass, our goal was to be liberal with paper inclusion, and to minimize false negatives. We interpreted program analysis tools in the broadest sense, to include model checkers, verifiers, static analysis tools, and dynamic tools. In the second pass, we examined the contents of the paper to identify if the paper did contribute or discuss its output for programmers. Finally, we removed papers that were purely related to reducing false positives, unless those papers used false positives as part of their output to provide additional information to a programmer. For some papers, the output was measured in terms of manual patches submitted to bug repositories. We excluded such papers since the output was manually constructed, and not directly from the tool.

Synthesis and Analysis. We synthesized the papers into a taxonomy of presentation attributes (Section 3). For analysis, we opted for a narrative approach in which we summarized the contributions of each of the papers with respect to human-friendly presentations.

2.3 Limitations

As a form of interim guidance, a scoping review has several important limitations, which we openly acknowledge. First, the review is biased in several ways. Being scoped only to PLDI means that the identified taxonomy is likely to be incomplete. Second, the scoping review by definition misses key contributions found in other conferences, such as the International Conference on Software Engineering (ICSE), Foundations of Software Engineering (FSE), and the Conference on Human Factors in Computing Systems (CHI), just to name a few. Third, the paper summaries are intended to be accessible to HCI researchers who may not have formal PL experience. As a result, in the interest of being broadly accessible, some of the summaries of the papers may be oversimplified in terms of their PL contributions. Finally, any conclusions made from this interim work should be treated as provisional and subject to revision as more comprehensive reviews are conducted.

3 Taxonomy of Presentation

In this section, we classify each of the identified papers that discuss or contribute to program analysis output intended for programmers. Intentionally, we labeled the taxonomy features such that they do not commit to a particular textual or visualization affordance. For example, in a text-interface, the feature of ranking (Section 3.0.7) may be implemented as an enumerated list of items in the console, with a prompt for selection, if interactivity is required. In a graphical interface, ranking might instead be implemented through a pop-up

Table 1. Taxonomy of Presentation

Feature	Section
Alignment	Section 3.0.1
Clustering and Classification	Section 3.0.2
Comparing	Section 3.0.3
Example	Section 3.0.4
Interactivity	Section 3.0.5
Localizing	Section 3.0.6
Ranking	Section 3.0.7
Reducing	Section 3.0.8
Tracing	Section 3.0.9

drop-down, through which the programmer would select their desired option.

The identified taxonomy of presentation features are summarized in Table 1. Some papers describe output that use multiple features; in such a case, we selected the feature which we felt best represented the contribution of the output.

3.0.1 Alignment

In alignment, program analysis output is presented in a representation that is already familiar to the programmer.

Within this feature, Pombrio and Krishnamurthi tackle the problem of syntactic sugar—programming constructs that make things easier to express—but are ultimately reducible to alternative constructs. For example, in C, the array access notation $a[i]$ is syntactic sugar for (the sometimes less convenient notation) $*(a + i)$. Unfortunately, syntactic sugar is eliminated by many transformation algorithms, making the resulting program unfamiliar to the programmer. Pombrio and Krishnamurthi introduce a process of *resugaring* to allow computation reductions in terms of the surface syntax [46]. With similar aims, the AutoCorres tool uses a technique of *specification abstraction*, to present programmers with a representation of the program at a human-readable abstraction while additionally producing a formal refinement of the final presentation [22].

Notions of natural language and readability find their place in several PLDI papers. Qiu et al. propose *natural proofs*, in which automated reasoning systems restrict themselves to using common patterns found in human proofs [44, 47]. Given a reference implementation, and an error model of potential corrections, Singh et al. introduce a method for automatically deriving *minimal corrections* to students’ incorrect solutions, in the form a itemized list of changes, expressed in natural language form [53]. And the AFix tool uses a variety of static analysis and static code transformations to design bug fixes for a type of concurrency bug, *single-variable atomicity violations* [28]. The bug fixes are human-friendly in that they attempt to provide a fix that, in addition to other

metrics, does not harm code *readability*. To support readability, the authors manually evaluated several possible locking policies to determine which ones were most readable.

Issues of alignment and representation become important to programmers during understanding of optimizations in source-level debugging of optimized code [2]; in their approach, Adl-Tabatabai and Gross implement *engendered variables* that would cause the programmer to draw incorrect conclusions as a result of internal optimizations by the compiler. Earlier work by [10] [10] and Coutant et al. [16] also provide techniques within this space.

3.0.2 Clustering and Classification

Clustering and classification output aims to organize or separate information in a way that reduces the cognitive burden for programmers. For example, Narayanasamy et al. focus on a dynamic analysis technique to automatically classify *data races*—a type of concurrency bug in multi-threaded programs—as being potentially benign or potentially harmful [40]; furthermore, the tool provides the programmer with a reproducible scenario of the data race to help the developer understand how it manifests.

Liblit et al. present a statistical dynamic debugging technique that *isolates* bugs in programs containing multiple undiagnosed bugs [36]; importantly, the algorithm separates the effects of different bugs and identifies predictors that are associated with individual bugs. An earlier technique using statistical sampling is also presented by the authors [35]. Other classification techniques include Ha et al.; they introduce CLARIFY, a system which classifies *behavior profiles*—essentially, an application’s behavior—for black box software components where the source code is not available [23]. And Ammons et al. consider the problem of specification on programs in that the specifications themselves need methods for debugging; they present a method for debugging formal, temporal specifications through *concept analysis* to automatically group traces into highly similar clusters [3].

3.0.3 Comparing

Comparisons occur in program analysis tools when the programmer has a need to examine or understand differences between two or more versions of their code. Within this feature, Hoffman et al. introduce a technique of *semantic views* of program executions to perform trace analysis; they apply their technique to identify regressions in large software applications [25]. Through a differencing technique, their RPRISM tool outputs a semantic “diff” between the original and new versions, to allow potential causes to be viewed in their full context. Similarly, early work by Horwitz identifies both semantic and textual differences between two versions of a program [26], in contrast to traditional diff-tools that treat source as plain text.

3.0.4 Example

Examples and counterexamples are forms of output that provide evidence for why a situation can occur or how a situation can be violated. Examples are usually provided in conjunction with other presentation features. Contributions in this feature focus on the type of example to present to the programmer, which is sometimes arbitrary and sometimes based on an output measure, such as minimizing lines of code.

The Alive-Infer tool, for example, infers preconditions to ensure the validity of a peephole compiler optimization. To the user, it reports both a *weakest* precondition and a set of “succinct” partial preconditions. For wrong optimizations, the tool provides counterexamples [39]. Zhang et al. apply a technique of *skeletal program enumeration* to generate small test programs for reporting bugs about in GCC and clang compilers; the generated test programs contain fewer than 30 lines on average [61]. Still other work with test programs devise a test-case reducer for C compiler bugs to obtain small and valid test-cases consistently [48]; the underlying machinery is based on *generic fixpoint computations* which invokes a *modular reducer*.

Padon et al. hypothesize that one of the reasons automated methods are difficult to use in practice is because they are opaque. As Padon et al. states, “they fail in ways that are difficult for a human user to understand and to remedy” [43]. Their system, Ivy, graphically displays concrete counterexamples to induction, and allows the user to interactively guide generation from these counterexamples [43]. Nguy  n and Van Horn implement a tool in Racket to generate counterexamples for erroneous modules [41] and Isradisaikul and Myers design an algorithm that generates helpful counterexamples for parsing ambiguities; for every parsing conflict, the algorithm generates a compact counterexample illustrating the ambiguity [27].

PSKETCH is a program synthesis tool that helps programmers implement concurrent data structures; it uses a *counter example guided inductive synthesis algorithm* (CEGIS) to converge to a solution within a handful of iterations [54].

For type error messages, Lerner et al. pursue an approach in which the type-checker itself does not produce error messages, but instead relies on an oracle for a search procedure that finds similar programs that *do* type-check; to bypass the typically-inscrutable type error messages, their system provides examples of code (at the same location) that would type check [33].

And for memory-related output, Cherem et al. implement a practical analysis algorithm for detecting memory leaks in C programs; their analysis uses *sparse value-flows* to present *concise* error messages to developers [13].

3.0.5 Interactivity

We identified several papers whose tools support interactivity in limited ways. That is, the programmer can interact with the tool either before the output is produced, in order to customize the output—or work with the output of the tool in a *mixed-initiative* fashion, where both the programmer and the tool collaborate to arrive at a solution.

Within this feature, Parsify is a program synthesis tool that synthesizes a parser from input and output examples. The tool interface provides immediate visual feedback in response to changes in the grammar being refined, as well as a graphical mechanism for specifying example parse trees using only textual selections [34].

Live programming is a user interface capability that allows a programmer to edit code and immediately see the effect of the code changes. Burckhardt et al. introduce a type and effect system that formalizes the separation of *rendering* and *non-rendering aspects* of the user interface to make feedback responsive [11].

Dillig et al. present a technique called *abductive inference*—that is, to find an explanatory hypothesis for a desired outcome—to assist programmers in classifying error reports. The technique computes small relevant queries presented to a user that capture exactly the information the analysis is missing to either discharge or validate the error [18].

LeakChaser identifies unnecessarily-held memory references which often result in memory leaks and performance issues in manages languages such as Java. The tool allows an *iterative* process through three *tiers* which assist programmers at different levels of abstraction, from *transactions* at the highest-level tier to *lifetime relationships* at the lowest level tier.

CHAMELEON assists programmers in choosing an abstract collection implementation in their algorithm [52]. During program execution, CHAMELEON computes trace metrics using *semantic profiling*, together with a set of collection selection rules, to present recommended collection adaptation strategies to the programmers. Similarly, the PetaBricks tool makes algorithm choice a first-class construct of the language [4].

von Dincklage and Diwan identify that too many underlying false positives in tools, such as in refactoring, can compromise a tool’s usefulness [59]. They propose a method to produce necessary and sufficient reasons, that is, a *why* explanation, for a potentially undesirable result; the programmer can then—through applying predicates—provide feedback on whether the given analysis result is desirable.

Finally, MrSpidey is a user-friendly, static debugger for Scheme [20]; the program analysis computes *value set descriptions* for each term in the program and constructs a *value flow graph* connecting the set descriptions; these flows are made visible to the programmer through a value flow browser which overlays arrows over the program text. The

441 programmer can interactively expose portions of the value
442 graph.

443 3.0.6 Localizing

445 We identified localizing in two forms: 1) a *point* localization,
446 in which a program analysis tool tries to identify a single
447 source as relevant to the error, and 2) as *slices*, where multiple
448 statements are responsible for the error.

449 **Point.** Zhang et al. implement, within the GHC compi-
450 ller, a simple Bayesian type error diagnostic that identifies
451 the *most likely* source of the type error, rather than the *first*
452 *source* the inference engine “trips over” [60]. The BugAssist
453 tool implements an algorithm for error cause localization
454 based on a reduction to *the maximal satisfiability problem*
455 to identify the *cause* of an error from failing execution tra-
456 ces [31]. The Breadcrumbs tool uses a *probabilistic calling*
457 *context* (essentially, a stack trace) to identify the root cause
458 of bug, by recording extra information that *might* be useful
459 in explaining a failure [9].

460 **Slices.** Program slicing identifies parts of the program
461 that may affect a point of interest—such as those related
462 to an error message; Sridharan et al. propose a technique
463 called *thin slicing* which helps programmer better identify
464 bugs because it identifies more relevant lines of code than
465 traditional slicing [55]. Analogous to thin slicing, Zhang
466 et al. developed a strategy for pruning dynamic slices to
467 identify subsets of statements that are likely responsible for
468 producing an incorrect value; for each statement executed
469 in the dynamic slice, their tool computes a confidence value,
470 with higher values corresponding to greater likelihood that
471 the execution of the statement produced a correct value [62].

472 3.0.7 Ranking

473 Ranking is a presentation feature that orders the output
474 of the program analysis in a systematic way. For example,
475 random testing tools, that is, *fuzzers*, can be frustrating to
476 use because they “indiscriminately and repeatedly find bugs
477 that may not be severe enough to fix right away” [12]. Chen
478 et al. propose a technique that *orders* test cases in a way that
479 diverse, interesting cases (defined through a machine techni-
480 que called *furthest point first*) are highly ranked [12]. And
481 the AcSPEC tool prioritizes alarms for automatic program
482 verifiers through *semantic inconsistency detection* in order to
483 report high-confidence warnings to the programmer [8].

484 Coppa et al. present a profiling methodology and toolkit
485 for helping developers discover asymptotic inefficiencies in
486 their code [15]. The output of the profiler is, for each exe-
487 cuted routine of the program, a set of tuples that aggregate
488 performance costs by input size—these outputs are intended
489 to be used as input to performance plots. The Kremlin tool
490 makes recommendations about which parts of the program a
491 programmer should spend effort parallelizing; the tool identi-
492 fies these regions through a *hierarchical critical path analysis*

496 and presents to the programmer an ordered (by speedup)
497 parallelism plan as a list of files and lines to modify [21].

498 Perelman et al. provide ranked expressions for completi-
499 ons in API libraries through a language of *partial expressions*,
500 which allows the programmer to leave “holes” for the parts
501 they do not know [45].

502 3.0.8 Reduction

503 Reduction approaches take a large design space of allowable
504 program output and reduce that space using some syste-
505 matic rule. For example, Logozzo et al. introduce a static
506 analysis technique of *Verification Modulo Versions* (VMV),
507 which reduces the number of alarms reported by verifiers
508 while maintaining semantic guarantees [38]. Specifically,
509 VMV is designed for scenarios in which programmers desire
510 to fix *new* defects introduced since a previous release.

511 3.0.9 Tracing

512 Tracing involves flows of information, and understanding
513 how information propagates across source code. As one ex-
514 ample, Ohmann et al. present a system that answers control-
515 flow queries posed by developers as formal languages. The
516 tool indicates whether the query expresses control flow that
517 is *possible* or *impossible* for a given failure report. As another
518 example, PIDGIN is a program analysis and understanding
519 tool that allows programmers to interactively explore *infor-*
520 *mation flows*, through *program dependence graphs*, within
521 their applications and investigate counterexamples [29]. Taint
522 analysis is another information-flow analysis that establis-
523 hes whether values from unstructured parameters may flow
524 into security-sensitive operations [58]; implemented as TAJ,
525 the tool additionally eliminates redundant reports through
526 hybrid thin slicing and *remediation logic* over library local
527 points. Other techniques, such as those by Rubio-González
528 et al., use data-flow analysis techniques to track errors as
529 they propagate through file system code [49].

530 To support algorithmic debugging, Faddegon and Chitil
531 developed a library in Haskell, that, after annotating sus-
532 pected functions, presents a detailed *computational tree* [19].
533 Computational trees are essentially a trace to help developers
534 understand how a program works or why it does not work.
535 The tool TraceBack provides debugging information for pro-
536 duction systems by providing execution history data about
537 program problems [6]; it uses *first-fault diagnosis* to discover
538 what went wrong the *first* time the fault is encountered.

540 MemSAT helps programmers debug and reason about
541 memory models: given an *axiomatic* specification, the tool
542 outputs a *trace*—sequences of reads and writes—of the pro-
543 gram in which the specification is satisfied, or a *minimal*
544 subset of the memory model and program constraints that
545 are unsatisfiable [56].

546 The Merlin security analysis tool infers *information flows*
547 in a program to identify security vulnerabilities such as cross-
548 site scripting and SQL inject attacks [37]. Internally, the

inference is based on modeling a *data propagation graph* using *probabilistic constraints*.

4 Discussion

Lack of user evaluations in PL. Although we identified and classified papers within PLDI in terms of a taxonomy of presentation, our investigation confirms that papers either perform no usability evaluation with programmers, or the claims of usability of the tool are made through intuition, using the authors of the paper as subjects. For example, consider the presentation attribute of alignment (Section 3.0.1), in which several assumptions are made about how output should be presented in familiar representations to the programmer. All of these assumptions appear to be intuitive—give output in the same level of syntactic sugar as their source code for consistency, use proof constructions commonly found in human proofs, and support readability. Unfortunately, none of these assumptions are tested with actual developers, reminding us of the concerns noted by Hanenberg and others in the introduction. It's likely some of these assumptions *are* actually incorrect, which may explain the lack of adoption in practice and the confusing tool output programmers report for many of these sophisticated program analysis tools.

Lack of operational tools in HCI. At the same time, HCI researchers perform usability studies on user interfaces, yet the experiments they conduct are often performed on prototype platforms that are built specifically for the experiment under consideration. Consequently, even if the user interfaces are found to be effective or usable for some measures, the tools themselves cannot actually be used in practice. Regrettably, this means that user interface advances remain within academic papers, and do not ever make it to actual programmers without significant investment for tools that may not even be possible to build due to fundamental technical limitations.

Bridging PL and HCI. In our view, both deficiencies in HCI and PL can be reduced by fostering collaborations between the disciplines. A cross-disciplinary approach to tool development would enable usable program analysis tools, by having a pipeline from program analysis tools to user evaluations in HCI. HCI contributions could then feedback to PL to further improve the output of program analysis tools. But doing so requires a cross-disciplinary community that can provide such opportunities for collaboration. We suggest that PLATEAU has the potential to become this community.

5 Conclusions

In this paper, we conducted a scoping review of PLDI from the period of 1988-2017, to identify and catalog papers for program analysis tools that discussed or made contributions

to the presentation of output towards programmers. Admittedly, a scoping review is only a starting point for investigation, and can only provide interim guidance. Nevertheless, our hope is that the scoping review we have conducted can serve to bootstrap a future comprehensive systematic literature reviews. We are open to feedback on practical methods to realizing that goal.

References

- [1] Gregory D Abowd and Russell Beale. 1991. Users, systems and interfaces: A unifying framework for interaction. In *People and Computers VI*. 73–87.
- [2] Ali-Reza Adl-Tabatabai and Thomas Gross. 1996. Source-level Debugging of Scalar Optimized Code. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation (PLDI '96)*. ACM, 33–43. <https://doi.org/10.1145/231379.231388>
- [3] Glenn Ammons, David Mandelin, Rastislav Bodik, and James R. Larus. 2003. Debugging Temporal Specifications with Concept Analysis. In *Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation (PLDI '03)*. ACM, 182–195. <https://doi.org/10.1145/781131.781152>
- [4] Jason Ansel, Cy Chan, Yee Lok Wong, Marek Olszewski, Qin Zhao, Alan Edelman, and Saman Amarasinghe. 2009. PetaBricks: A Language and Compiler for Algorithmic Choice. In *Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '09)*. ACM, 38–49. <https://doi.org/10.1145/1542476.1542481>
- [5] H Arksey and L O'Malley. 2005. Scoping studies: towards a methodological framework. *Int J Soc Res Methodol* 8 (2005).
- [6] Andrew Ayers, Richard Schooler, Chris Metcalf, Anant Agarwal, Junghwan Rhee, and Emmett Witchel. 2005. TraceBack: First Fault Diagnosis by Reconstruction of Distributed Control Flow. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '05)*. ACM, 201–212. <https://doi.org/10.1145/1065010.1065035>
- [7] Titus Barik, Justin Smith, Kevin Lubick, Elisabeth Holmes, Jing Feng, Emerson Murphy-Hill, and Chris Parnin. 2017. Do Developers Read Compiler Error Messages?. In *Proceedings of the 39th International Conference on Software Engineering (ICSE '17)*. IEEE Press, Piscataway, NJ, USA, 575–585. <https://doi.org/10.1109/ICSE.2017.59>
- [8] Sam Blackshear and Shuvendu K. Lahiri. 2013. Almost-correct Specifications: A Modular Semantic Framework for Assigning Confidence to Warnings. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*. ACM, 209–218. <https://doi.org/10.1145/2491956.2462188>
- [9] Michael D. Bond, Graham Z. Baker, and Samuel Z. Guyer. 2010. Breadcrumbs: Efficient Context Sensitivity for Dynamic Bug Detection Analyses. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '10)*. ACM, 13–24. <https://doi.org/10.1145/1806596.1806599>
- [10] Gary Brooks, Gilbert J. Hansen, and Steve Simmons. 1992. A New Approach to Debugging Optimized Code. In *Proceedings of the ACM SIGPLAN 1992 Conference on Programming Language Design and Implementation (PLDI '92)*. ACM, 1–11. <https://doi.org/10.1145/143095.143108>
- [11] Sebastian Burckhardt, Manuel Fahndrich, Peli de Halleux, Sean McDermid, Michal Moskal, Nikolai Tillmann, and Jun Kato. 2013. It's Alive! Continuous Feedback in UI Programming. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*. ACM, 95–104. <https://doi.org/10.1145/2491956.2462170>
- [12] Yang Chen, Alex Groce, Chaoqiang Zhang, Weng-Keen Wong, Xiaoli Fern, Eric Eide, and John Regehr. 2013. Taming Compiler Fuzzers. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming*

- 661 *Language Design and Implementation (PLDI '13)*. ACM, 197–208. <https://doi.org/10.1145/2491956.2462173>
- 662
- 663 [13] Sigmund Cheren, Lonnie Princehouse, and Radu Rugina. 2007. Practical Memory Leak Detection Using Guarded Value-flow Analysis. In *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '07)*. ACM, 480–491. <https://doi.org/10.1145/1250734.1250789>
- 664
- 665
- 666
- 667 [14] Maria Christakis and Christian Bird. 2016. What developers want and need from program analysis: an empirical study. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering - ASE 2016*. ACM Press, New York, New York, USA, 332–343. <https://doi.org/10.1145/2970276.2970347>
- 668
- 669
- 670
- 671 [15] Emilio Coppa, Camil Demetrescu, and Irene Finocchi. 2012. Input-sensitive Profiling. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '12)*. ACM, 89–98. <https://doi.org/10.1145/2254064.2254076>
- 672
- 673
- 674 [16] D. S. Coutant, S. Meloy, and M. Russetta. 1988. DOC: A Practical Approach to Source-level Debugging of Globally Optimized Code. In *Proceedings of the ACM SIGPLAN 1988 Conference on Programming Language Design and Implementation (PLDI '88)*. ACM, 125–134. <https://doi.org/10.1145/53990.54003>
- 675
- 676
- 677 [17] Natasha Danas, Tim Nelson, Lane Harrison, Shriram Krishnamurthi, and Daniel J. Dougherty. 2017. User Studies of Principled Model Finder Output. In *Software Engineering and Formal Methods, (Software Engineering and Formal Methods,)*.
- 678
- 679
- 680
- 681
- 682 [18] Isil Dillig, Thomas Dillig, and Alex Aiken. 2012. Automat-d Error Diagnosis Using Abductive Inference. In *Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '12)*. ACM, 181–192. <https://doi.org/10.1145/2254064.2254087>
- 683
- 684
- 685
- 686 [19] Maarten Faddegon and Olaf Chitil. 2016. Lightweight Computation Tree Tracing for Lazy Functional Languages. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '16)*. ACM, 114–128. <https://doi.org/10.1145/2908080.2908104>
- 687
- 688
- 689
- 690 [20] Cormac Flanagan, Matthew Flatt, Shriram Krishnamurthi, Stephanie Weirich, and Matthias Felleisen. 1996. Catching Bugs in the Web of Program Invariants. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation (PLDI '96)*. ACM, 23–32. <https://doi.org/10.1145/231379.231387>
- 691
- 692
- 693
- 694 [21] Saturnino Garcia, Donghwan Jeon, Christopher M. Louie, and Michael Bedford Taylor. 2011. Kremlin: Rethinking and Rebooting Gprof for the Multicore Age. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '11)*. ACM, 458–469. <https://doi.org/10.1145/1993498.1993553>
- 695
- 696
- 697
- 698 [22] David Greenaway, Japheth Lim, June Andronick, and Gerwin Klein. 2014. Don'T Sweat the Small Stuff: Formal Verification of C Code Without the Pain. *SIGPLAN Not.* 49, 6 (June 2014), 429–439. <https://doi.org/10.1145/2666356.2594296>
- 699
- 700
- 701
- 702 [23] Jungwoo Ha, Christopher J. Rossbach, Jason V. Davis, Indrajit Roy, Hany E. Ramadan, Donald E. Porter, David L. Chen, and Emmett Witchel. 2007. Improved Error Reporting for Software That Uses Black-box Components. In *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '07)*. ACM, 101–111. <https://doi.org/10.1145/1250734.1250747>
- 703
- 704
- 705
- 706 [24] Stefan Hanenberg. 2010. Faith, Hope, and Love: An Essay on Software Science's Neglect of Human Factors. In *Proceedings of the ACM International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA '10)*. ACM, 933–946. <https://doi.org/10.1145/1869459.1869536>
- 707
- 708
- 709
- 710
- 711 [25] Kevin J. Hoffman, Patrick Eugster, and Suresh Jagannathan. 2009. Semantics-aware Trace Analysis. In *Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '09)*. ACM, 453–464. <https://doi.org/10.1145/1542476.1542527>
- 712
- 713
- 714
- 715
- 716 [26] Susan Horwitz. 1990. Identifying the Semantic and Textual Differences Between Two Versions of a Program. In *Proceedings of the ACM SIGPLAN 1990 Conference on Programming Language Design and Implementation (PLDI '90)*. ACM, 234–245. <https://doi.org/10.1145/93542.93574>
- 717
- 718
- 719
- 720 [27] Chinawat Isradisaikul and Andrew C. Myers. 2015. Finding Counterexamples from Parsing Conflicts. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '15)*. ACM, 555–564. <https://doi.org/10.1145/2737924.2737961>
- 721
- 722
- 723
- 724 [28] Guoliang Jin, Linhai Song, Wei Zhang, Shan Lu, and Ben Liblit. 2011. Automated Atomicity-violation Fixing. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '11)*. ACM, 389–400. <https://doi.org/10.1145/1993498.1993544>
- 725
- 726
- 727
- 728 [29] Andrew Johnson, Lucas Wayne, Scott Moore, and Stephen Chong. 2015. Exploring and Enforcing Security Guarantees via Program Dependence Graphs. *SIGPLAN Not.* 50, 6 (June 2015), 291–302. <https://doi.org/10.1145/2813885.2737957>
- 729
- 730
- 731 [30] Brittany Johnson, Yoonki Song, Emerson Murphy-Hill, and Robert Bowdidge. 2013. Why don't software developers use static analysis tools to find bugs?. In *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 672–681. <https://doi.org/10.1109/ICSE.2013.6606613>
- 732
- 733
- 734
- 735 [31] Manu Jose and Rupak Majumdar. 2011. Cause Clue Clauses: Error Localization Using Maximum Satisfiability. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '11)*. ACM, 437–446. <https://doi.org/10.1145/1993498.1993550>
- 736
- 737
- 738
- 739 [32] Barbara Kitchenham. 2004. Procedures for performing systematic reviews. *Keele, UK, Keele University* 33, 2004 (2004), 1–26.
- 740
- 741 [33] Benjamin S. Lerner, Matthew Flower, Dan Grossman, and Craig Chambers. 2007. Searching for Type-error Messages. In *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '07)*. ACM, 425–434. <https://doi.org/10.1145/1250734.1250783>
- 742
- 743
- 744
- 745 [34] Alan Leung, John Sarracino, and Sorin Lerner. 2015. Interactive Parser Synthesis by Example. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '15)*. ACM, 565–574. <https://doi.org/10.1145/2737924.2738002>
- 746
- 747
- 748
- 749 [35] Ben Liblit, Alex Aiken, Alice X. Zheng, and Michael I. Jordan. 2003. Bug Isolation via Remote Program Sampling. In *Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation (PLDI '03)*. ACM, 141–154. <https://doi.org/10.1145/781131.781148>
- 750
- 751
- 752
- 753 [36] Ben Liblit, Mayur Naik, Alice X. Zheng, Alex Aiken, and Michael I. Jordan. 2005. Scalable Statistical Bug Isolation. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '05)*. ACM, 15–26. <https://doi.org/10.1145/1065010.1065014>
- 754
- 755
- 756
- 757 [37] Benjamin Livshits, Aditya V. Nori, Sriram K. Rajamani, and Anindya Banerjee. 2009. Merlin: Specification Inference for Explicit Information Flow Problems. In *Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '09)*. ACM, 75–86. <https://doi.org/10.1145/1542476.1542485>
- 758
- 759
- 760
- 761 [38] Francesco Logozzo, Shuvendu K. Lahiri, Manuel Fähndrich, and Sam Blackshear. 2014. Verification Modulo Versions: Towards Usable Verification. *SIGPLAN Not.* 49, 6 (June 2014), 294–304. <https://doi.org/10.1145/2666356.2594326>
- 762
- 763
- 764
- 765 [39] David Menendez and Santosh Nagarakatte. 2017. Alive-Infer: Data-driven Precondition Inference for Peephole Optimizations in LLVM. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2017)*. ACM, 49–63. <https://doi.org/10.1145/3062341.3062372>
- 766
- 767
- 768
- 769
- 770

- 771 [40] Satish Narayanasamy, Zhenghao Wang, Jordan Tigani, Andrew Ed- 826
772 wards, and Brad Calder. 2007. Automatically Classifying Benign and 827
773 Harmful Data Races Using Replay Analysis. In *Proceedings of the 28th* 828
774 *ACM SIGPLAN Conference on Programming Language Design and* 829
775 *Implementation (PLDI '07)*. ACM, 22–31. [https://doi.org/10.1145/1250734.](https://doi.org/10.1145/1250734.1250738)
776 1250738
- 776 [41] Phúc C Nguyễn and David Van Horn. 2015. Relatively Complete 830
777 Counterexamples for Higher-order Programs. *SIGPLAN Not.* 50, 6 831
778 (June 2015), 446–456. <https://doi.org/10.1145/2813885.2737971> 832
- 779 [42] Peter Ohmann, Alexander Brooks, Loris D'Antoni, and Ben 833
780 Liblit. 2017. Control-flow Recovery from Partial Failure Reports. In 834
781 *Proceedings of the 38th ACM SIGPLAN Conference on Programming* 835
782 *Language Design and Implementation (PLDI 2017)*. ACM, 390–405. <https://doi.org/10.1145/3062341.3062368> 836
- 783 [43] Oded Padon, Kenneth L. McMillan, Aurojit Panda, Mooly Sagiv, and 837
784 Sharon Shoham. 2016. Ivy: Safety Verification by Interactive General- 838
785 ization. In *Proceedings of the 37th ACM SIGPLAN Conference on* 839
786 *Programming Language Design and Implementation (PLDI '16)*. ACM, 614–630. <https://doi.org/10.1145/2908080.2908118> 840
- 787 [44] Edgar Pek, Xiaokang Qiu, and P. Madhusudan. 2014. Natural Proofs 841
788 for Data Structure Manipulation in C Using Separation Logic. *SIG-* 842
789 *PLAN Not.* 49, 6 (June 2014), 440–451. [https://doi.org/10.1145/2666356.](https://doi.org/10.1145/2666356.2594325)
790 2594325
- 791 [45] Daniel Perelman, Sumit Gulwani, Thomas Ball, and Dan Grossman. 843
792 2012. Type-directed Completion of Partial Expressions. In *Proceedings* 844
793 *of the 33rd ACM SIGPLAN Conference on Programming Language Design* 845
794 *and Implementation (PLDI '12)*. ACM, 275–286. [https://doi.org/10.1145/2254064.](https://doi.org/10.1145/2254064.2254098)
795 2254098
- 796 [46] Justin Pombrio and Shriram Krishnamurthi. 2014. Resugaring: Lifting 846
797 Evaluation Sequences Through Syntactic Sugar. *SIGPLAN Not.* 49, 6 847
798 (June 2014), 361–371. <https://doi.org/10.1145/2666356.2594319> 848
- 799 [47] Xiaokang Qiu, Pranav Garg, Andrei Ştefănescu, and Parthasarathy 849
800 Madhusudan. 2013. Natural Proofs for Structure, Data, and Separation. 850
801 In *Proceedings of the 34th ACM SIGPLAN Conference on Programming* 851
802 *Language Design and Implementation (PLDI '13)*. ACM, New York, NY, 852
803 USA, 231–242. <https://doi.org/10.1145/2491956.2462169> 853
- 804 [48] John Regehr, Yang Chen, Pascal Cuoq, Eric Eide, Chucky Ellison, and 854
805 Xuejun Yang. 2012. Test-case Reduction for C Compiler Bugs. In 855
806 *Proceedings of the 33rd ACM SIGPLAN Conference on Programming* 856
807 *Language Design and Implementation (PLDI '12)*. ACM, 335–346. [https://doi.org/10.1145/2254064.](https://doi.org/10.1145/2254064.2254104)
808 2254104
- 809 [49] Cindy Rubio-González, Haryadi S. Gunawi, Ben Liblit, Remzi H. Arpacı- 857
810 Dusseau, and Andrea C. Arpacı-Dusseau. 2009. Error Propagation 858
811 Analysis for File Systems. In *Proceedings of the 30th ACM SIGPLAN* 859
812 *Conference on Programming Language Design and Implementation (PLDI* 860
813 *'09)*. ACM, 270–280. [https://doi.org/10.1145/1542476.](https://doi.org/10.1145/1542476.1542506)
814 1542506
- 815 [50] Holger J Schünemann and Lorenzo Moja. 2015. Reviews: Rapid! Rapid! 861
816 Rapid! ... and systematic. *Systematic Reviews* 4, 1 (2015), 4. <https://doi.org/10.1186/2046-4053-4-4> 862
- 817 [51] Hyunmin Seo, Caitlin Sadowski, Sebastian Elbaum, Edward Aftandilian, 863
818 and Robert Bowdidge. 2014. Programmers' build errors: a case 864
819 study (at google). In *Proceedings of the 36th International Conference* 865
820 *on Software Engineering - ICSE 2014*. ACM Press, New York, New York, 866
821 USA, 724–734. <https://doi.org/10.1145/2568225.2568255> 867
- 822 [52] Ohad Shacham, Martin Vechev, and Eran Yahav. 2009. Chameleon: 868
823 Adaptive Selection of Collections. In *Proceedings of the 30th ACM* 869
824 *SIGPLAN Conference on Programming Language Design and Imple-* 870
825 *mentation (PLDI '09)*. ACM, 408–418. [https://doi.org/10.1145/1542476.](https://doi.org/10.1145/1542476.1542522)
826 1542522
- 826 [53] Rishabh Singh, Sumit Gulwani, and Armando Solar-Lezama. 2013. 871
827 Automated Feedback Generation for Introductory Programming As- 872
828 signments. *SIGPLAN Not.* 48, 6 (June 2013), 15–26. [https://doi.org/10.1145/2499370.](https://doi.org/10.1145/2499370.2462195)
829 2462195
- 830 [54] Armando Solar-Lezama, Christopher Grant Jones, and Rastislav Bodik. 873
831 2008. Sketching Concurrent Data Structures. In *Proceedings of the* 874
832 *29th ACM SIGPLAN Conference on Programming Language Design and* 875
833 *Implementation (PLDI '08)*. ACM, 136–148. [https://doi.org/10.1145/1375581.](https://doi.org/10.1145/1375581.1375599)
834 1375599
- 834 [55] Manu Sridharan, Stephen J. Fink, and Rastislav Bodik. 2007. Thin Slic- 876
835 ing. In *Proceedings of the 28th ACM SIGPLAN Conference on Program-* 877
836 *ming Language Design and Implementation (PLDI '07)*. ACM, 112–122. 878
837 <https://doi.org/10.1145/1250734.1250748> 879
- 838 [56] Emina Torlak, Mandana Vaziri, and Julian Dolby. 2010. MemSAT: 880
839 Checking Axiomatic Specifications of Memory Models. *SIGPLAN Not.* 45, 6 (June 2010), 341–350. <https://doi.org/10.1145/1809028.1806635>
- 840 [57] V. Javier Traver. 2010. On compiler error messages: What they say 881
841 and what they mean. *Advances in Human-Computer Interaction* 2010 882
842 (2010), 1–26. <https://doi.org/10.1155/2010/602570>
- 843 [58] Omer Tripp, Marco Pistoia, Stephen J. Fink, Manu Sridharan, and Omri 883
844 Weisman. 2009. TAJ: Effective Taint Analysis of Web Applications. 884
845 In *Proceedings of the 30th ACM SIGPLAN Conference on Programming* 885
846 *Language Design and Implementation (PLDI '09)*. ACM, 87–97. [https://doi.org/10.1145/1542476.](https://doi.org/10.1145/1542476.1542486)
847 1542486
- 847 [59] Daniel von Dincklage and Amer Diwan. 2008. Explaining Failures of 886
848 Program Analyses. In *Proceedings of the 29th ACM SIGPLAN Conference* 887
849 *on Programming Language Design and Implementation (PLDI '08)*. ACM, 260–269. <https://doi.org/10.1145/1375581.1375614> 888
- 850 [60] Danfeng Zhang, Andrew C. Myers, Dimitrios Vytiniotis, and Simon 889
851 Peyton-Jones. 2015. Diagnosing Type Errors with Class. In *Proceedings* 890
852 *of the 36th ACM SIGPLAN Conference on Programming Language Design* 891
853 *and Implementation (PLDI '15)*. ACM, 12–21. [https://doi.org/10.1145/2737924.](https://doi.org/10.1145/2737924.2738009)
854 2738009
- 854 [61] Qirun Zhang, Chengnian Sun, and Zhendong Su. 2017. Skeletal Pro- 892
855 gram Enumeration for Rigorous Compiler Testing. In *Proceedings of* 893
856 *the 38th ACM SIGPLAN Conference on Programming Language Design* 894
857 *and Implementation (PLDI 2017)*. ACM, 347–361. <https://doi.org/10.1145/3062341.3062379> 895
- 858 [62] Xiangyu Zhang, Neelam Gupta, and Rajiv Gupta. 2006. Pruning Dynam- 896
859 ic Slices with Confidence. In *Proceedings of the 27th ACM SIGPLAN* 897
860 *Conference on Programming Language Design and Implementation* 898
861 *(PLDI '06)*. ACM, 169–180. <https://doi.org/10.1145/1133981.1134002> 899