# Finding Bugs in Web Applications Using Dynamic Test Generation and Explicit-State Model Checking

Shay Artzi, Adam Kieżun, Julian Dolby, Frank Tip, Danny Dig,
Amit Paradkar, *Senior Member*, *IEEE*, and Michael D. Ernst

**Abstract**—Web script crashes and malformed dynamically generated webpages are common errors, and they seriously impact the usability of Web applications. Current tools for webpage validation cannot handle the dynamically generated pages that are ubiquitous on today's Internet. We present a dynamic test generation technique for the domain of dynamic Web applications. The technique utilizes both combined concrete and symbolic execution and explicit-state model checking. The technique generates tests automatically, runs the tests capturing logical constraints on inputs, and minimizes the conditions on the inputs to failing tests so that the resulting bug reports are small and useful in finding and fixing the underlying faults. Our tool Apollo implements the technique for the PHP programming language. Apollo generates test inputs for a Web application, monitors the application for crashes, and validates that the output conforms to the HTML specification. This paper presents Apollo's algorithms and implementation, and an experimental evaluation that revealed 673 faults in six PHP Web applications.

**Index Terms**—Software testing, Web applications, dynamic analysis, PHP, reliability, verification.

✦

---

## 1 INTRODUCTION

DYNAMIC test generation tools, such as DART [17], Cute [39], and EXE [7], generate tests by executing an application on concrete input values, and then creating additional input values by solving symbolic constraints derived from exercised control-flow paths. To date, such approaches have not been practical in the domain of Web applications, which pose special challenges due to the dynamism of the programming languages, the use of implicit input parameters, their use of persistent state, and their complex patterns of user interaction.

This paper extends dynamic test generation to the domain of web applications that dynamically create web (HTML) pages during execution, which are typically presented to the user in a browser. Apollo applies these techniques in the context of the scripting language PHP, one of the most popular languages for server-side Web programming. According to the Internet research service, Netcraft,[1] PHP powered 21 million domains as of April 2007, including

large, well-known websites such as Wikipedia and WordPress. In addition to dynamic content, modern Web applications may also generate significant application logic, typically in the form of JavaScript code that is executed on the client side. Our techniques are primarily focused on server-side PHP code, although we do some minimal analysis of client-side code to determine how it invokes additional server code through user-interface mechanisms such as forms.

Our goal is to find two kinds of failures in web applications: *execution failures* that are manifested as crashes or warnings during program execution, and *HTML failures* that occur when the application generates malformed HTML. Execution failures may occur, for example, when a web application calls an undefined function or reads a nonexistent file. In such cases, the HTML output contains an error message and execution of the application may be halted, depending on the severity of the failure. HTML failures occur when output is generated that is not syntactically well-formed HTML (e.g., when an opening tag is not accompanied by a matching closing tag). HTML failures are generally not as important as execution failures because Web browsers are designed to tolerate some degree of malformedness in HTML, but they are undesirable for several reasons. First and most serious is that browsers' attempts to compensate for malformed webpages may lead to crashes and security vulnerabilities.[2] Second, standard HTML renders faster.[3] Third, malformed HTML is less portable across browsers and is vulnerable to

breaking or looking strange when displayed by browser versions on which it is not tested. Fourth, a browser might succeed in displaying only part of a malformed webpage, while silently discarding important information. Fifth, search engines may have trouble indexing malformed pages [48].

Web developers widely recognize the importance of creating legal HTML. Many websites are checked using HTML validators.[4] However, HTML validators can only point out problems in HTML pages, and are by themselves incapable of finding faults in applications that *generate* HTML pages. Checking *dynamic* Web applications (i.e., applications that generate pages during execution) requires checking that the application creates a valid HTML page on *every* possible execution path. In practice, even professionally developed and thoroughly tested applications often contain multiple faults (see Section 6).

There are two general approaches to finding faults in web applications: static analysis and dynamic analysis (testing). In the context of Web applications, static approaches have limited potential because 1) Web applications are often written in dynamic scripting languages that enable on-the-fly creation of code, and 2) control in a Web application typically flows via the generated HTML text (e.g., buttons and menus that require user interaction to execute), rather than solely via the analyzed code. Both of these issues pose significant challenges to approaches based on static analysis. Testing of dynamic Web applications is also challenging because the input space is large and applications typically require multiple user interactions. The state of the practice in validation for Web-standard compliance of real Web applications involves the use of programs such as HTML Kit[5] that validate each generated page, but require manual generation of inputs that lead to displaying different pages. We know of no automated tool that automatically generates inputs that exercise different control-flow paths in a Web application, and validates the dynamically generated HTML pages that the Web application generates when those paths are executed.

This paper presents an automated technique for finding failures in HTML-generating web applications. Our technique is based on dynamic test generation, using combined concrete and symbolic (concolic) execution, and constraint solving [7], [17], [39]. We created a tool, Apollo, that implements our technique in the context of the publicly available PHP interpreter.

Apollo first executes the Web application under test with an empty input. During each execution, Apollo monitors the program to record *path constraints* that reflect how input values affect control flow. Additionally, for each execution, Apollo determines whether execution failures or HTML failures occur (for HTML failures, an HTML validator is used as an oracle). Apollo automatically and iteratively creates new inputs using the recorded path constraints to create inputs that exercise different control flow. Most previous approaches for concolic execution only detect "standard errors" such as crashes and assertion failures. Our approach detects such standard errors as well, but also uses an oracle to detect specification violations in the application's output.

Another novelty in our work is the inference of input parameters, which are not manifested in the source code,

but which are interactively supplied by the user (e.g., by clicking buttons in generated HTML pages). The desired behavior of a PHP application is usually achieved by a series of interactions between the user and the server (e.g., a minimum of five user actions are needed from opening the main Amazon page to buying a book). We handle this problem by enhancing the combined concrete and symbolic execution technique with explicit-state model checking based on automatic dynamic simulation of user interactions. In order to simulate user interaction, Apollo stores the state of the environment (database, sessions, and cookies) after each execution, analyzes the output of the execution to detect the possible user options that are available, and restores the environment state before executing a new script based on a detected user option.

Techniques based on combined concrete and symbolic executions [7], [17], [39] may create multiple inputs that expose the same fault. In contrast to previous techniques, to avoid overwhelming the developer, our technique automatically identifies the minimal part of the input that is responsible for triggering the failure. This step is similar in spirit to Delta Debugging [9]. However, since Delta Debugging is a general, *black box* input minimization technique, it is oblivious to the properties of inputs. In contrast, our technique is *white box*: It uses the information that certain inputs induce partially overlapping control-flow paths. By intersecting these paths, our technique significantly minimizes the constraints on the inputs.

The contributions of this paper are the following:

- We adapt the established technique of dynamic test generation, based on combined concrete and symbolic execution [7], [17], [39], to the domain of PHP Web applications. This involves:

  1. using an HTML verifier as an oracle to find errors in *dynamically generated* HTML,
  2. dynamically discovering possible input parameters,[6]
  3. dealing with data types and operations specific to the PHP language,
  4. tracking the use of persistent state and how input flows through it, and
  5. automatically discovering input values based on the examination of branch conditions on execution paths.

- We created a tool, Apollo, that implements the technique for PHP.
- We evaluated our tool by applying it to six real Web applications and comparing the results with random testing. We show that dynamic test generation can be effective when adapted to the domain of Web applications written in PHP: Apollo identified 673 faults while achieving line coverage of 52.9 percent.
- We present a detailed classification of the faults found by Apollo.

The remainder of this paper is organized as follows: Section 2 presents an overview of PHP, introduces our

---

4. http://validator.w3.org, http://www.htmlhelp.com/tools/validator.
5. http://www.htmlkit.com.

6. Halfond and Orso [20] use static analysis techniques to solve a similar problem in the context of Web applications written in Java.

```php
1  <?php
2
3  make_header(); // print HTML header
4
5  // Make the $page variable easy to use //
6  if(!isset($_GET['page'])) $page = 0;
7  else $page = $_GET['page'];
8
9  // Bring up the report cards and stop processing //
10 if($_GET['page2']==1337) {
11   require('printReportCards.php');
12   die();  // terminate the PHP program
13 }
14
15 // Validate and log the user into the system //
16 if($_GET["login"] == 1) validateLogin();
17
18 switch ($page)
19 {
20   case 0:  require('login.php'); break;
21   case 1:  require('TeacherMain.php'); break;
22   case 2:  require('StudentMain.php'); break;
23   default: die("Incorrect page number.  Please verify.");
24 }
25
26 make_footer(); // print HTML footer
27 ...
```

```php
27 function validateLogin() {
28   if(!isset($_GET['username'])) {
29     echo "<j2> username must be supplied.</h2>\n";
30     return;
31   }
32   $username = $_GET['username'];
33   $password = $_GET['password'];
34   if($username=="john" && $password=="theTeacher")
35     $page=1;
36   else if($username=="john" && $password=="theStudent")
37     $page=2;
38   else echo "<h2>Login error. Please try again</h2>\n";
39 }
40
41 function make_header() { // print HTML header
42 print("
43 <!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN"
44    "http://www.w3.org/TR/html4/strict.dtd">
45 <HTML>
46  <HEAD> <TITLE> Class Management  </TITLE> </HEAD>
47  <BODY>");
48 }
49
50 function make_footer() {  // close HTML elements opened by header()
51 print("
52  </BODY>
53 </HTML>");
54 }
55 ?>
```

Fig. 1. A simplified PHP program excerpt from SchoolMate. This excerpt contains three faults (two real, one seeded), explained in Section 2.3.

running example, and discusses classes of failures in PHP web applications. Section 3 presents a simplified version of the algorithm and illustrates it on an example program. Section 4 presents the complete algorithm handling stateful execution with the simulation of interactive user inputs and illustrates it on an example program. Section 5 discusses the implementation of Apollo. Section 6 presents our experimental evaluation of Apollo on open-source Web applications. Section 7 gives an overview of related work, and Section 8 presents conclusions.

## 2  CONTEXT: PHP WEB APPLICATIONS

### 2.1  The PHP Scripting Language

This section briefly reviews the PHP scripting language, focusing on those aspects of PHP that differ from mainstream languages. Readers familiar with PHP may skip to the discussion of the running example in Section 2.2.

PHP is widely used for implementing Web applications, in part due to its rich library support for network interaction, HTTP processing, and database access. The input to a PHP program is a map from strings to strings. Each key is a parameter that the program can read, write, or check if it is set. The string value corresponding to a key may be interpreted as a numerical value if appropriate. The output of a PHP Web application is an HTML document that can be presented in a Web browser.

PHP is object oriented, in the sense that it has classes, interfaces, and dynamically dispatched methods with syntax and semantics similar to that of Java. PHP also has features of scripting languages, such as dynamic typing and an `eval` construct that interprets and executes a string value that was computed at runtime as a code fragment. For example, the following code fragment:

$$\$code = `\$x = 3;';\ \$x = 7; eval(\$code); echo\ \$x;$$

prints the value 3 (names of PHP variables start with the $ character). Other examples of the dynamic nature of

PHP are a predicate that checks whether a variable has been defined, and class and function definitions are statements that may occur anywhere.

The code in Fig. 1 illustrates the flavor of PHP. The require statement that is used in line 11 of Fig. 1 resembles the C #include directive in the sense that it includes the code from another source file. However, the C version is a preprocessor directive with a constant argument, whereas the PHP version is an ordinary statement in which the file name is computed at runtime. There are many similar cases where runtime values are used, e.g., switch labels need not be constant. This degree of flexibility is prized by PHP developers for enabling rapid application prototyping and development. However, the flexibility can make the overall structure of program hard to discern and it can make programs prone to code quality problems.

### 2.2  PHP Example

The PHP program of Fig. 1 is a simplified version of SchoolMate,[7] which allows school administrators to manage classes and users, teachers to manage assignments and grades, and students to access their information.

Lines 6 and 7 read the global parameter page that is supplied to the program in the URL, e.g., http://www.mywebsite.com/index.php?page=1. Line 10 examines the value of the global parameter page2 to determine whether to evaluate file printReportCards.php.

Function validateLogin (lines 27-39) sets the global parameter page to the correct value based on the identity of the user. This value is used in the switch statement on line 18, which presents the login screen or one of the teacher/student screens.

### 2.3  Failures in PHP Programs

Our technique targets two types of failures that can be automatically identified during the execution of PHP web applications. First, *execution failures* may be caused by a

---
7. http://sourceforge.net/projects/schoolmate.

missing included file, an incorrect MySQL query, or an uncaught exception. Such failures are easily identified as the PHP interpreter generates an error message and halts execution. Less serious execution failures, such as those caused by the use of deprecated language constructs, produce obtrusive error messages but do not halt execution. The last kind of execution failures are execution halting calls signaling unclean exits, such as `die` with a nonempty message or `exit` with a nonzero value. Second, *HTML failures* involve situations in which the generated HTML page is not syntactically correct according to an HTML validator. Section 1 discussed several negative consequences of malformed HTML.

Since HTML failures might be the result of an earlier execution failure, Apollo ignores HTML failures in cases where the execution contains any execution failure.

As an example, the program of Fig. 1 contains three faults, which cause the following failures when the program is executed:

1. Executing line 11 of the program results in an *execution failure* because the file `printReport-Cards.php` referenced on that line is missing.
2. In some cases, the program execution may result in an execution failure when the default case of the switch statement on line 23 is executed. Specifically, line 23 terminates program execution uncleanly when the global parameter `page` is not 0, 1, or 2 and when `page` is not written by function `Validate Login`. The output in this case is *malformed HTML* because the `make_footer` method is not executed, resulting in an unclosed HTML tag in the output. However, this HTML failure is not reported since Apollo does not check for malformed HTML when the execution results in an execution failure.
3. The program produces *malformed HTML* when line 29 generates an illegal HTML tag `j2`.

The first failure is similar to a failure that our tool found in one of the PHP applications we studied. The second failure is caused by a fault that exists in the original code of the SchoolMate program. The third failure is the result of a fault that was artificially inserted into the example for illustration.

## 2.4 Completeness

The analysis presented in this paper is a dynamic analysis. There is no guarantee of completeness (in terms of the number of bugs found or in terms of code coverage) for several reasons, including the fact that the application being analyzed may contain dead code and because the tool may be unable to cover certain parts of the application. This may, for example, occur when the solver is incapable of inferring the conditions under which a given statement is reachable, and also because our analysis of user options (e.g., in JavaScript that is embedded in the generated HTML pages) is incomplete.

## 3 FINDING FAILURES IN PHP WEB APPLICATIONS

Our technique for finding failures in PHP applications is a variation on an established dynamic test generation technique [7], [17], [18], [39] sometimes referred to as concolic testing. For expository purposes, we will present the algorithm in two steps. First, this section presents a simplified version of the algorithm that does not simulate user inputs or keep track of persistent session state. We will demonstrate this simplified algorithm on the example of Fig. 1. Then, Section 4 presents a generalized version of the algorithm that handles user-input simulation and stateful executions, and illustrates it on a more complex example.

The basic idea behind the technique is to execute an application on some initial input (e.g., an arbitrarily or randomly chosen input), and then on additional inputs obtained by solving constraints derived from exercised control-flow paths. We adapted this technique to PHP Web applications as follows:

- We extend the technique to consider failures other than execution failures by using an oracle to determine whether or not program output is correct. In particular, we use an HTML validator to determine whether the output is a well-formed HTML page.
- The PHP language contains constructs such as `isset` (checking whether a variable is defined), `isempty` (checking whether a variable contains a value from a specific set), `require` (dynamic loading of additional code to be executed), `header` for redirection of execution, and several others that require the generation of constraints that are absent in languages such as C or Java.
- PHP applications typically interact with a database and need appropriate values for user authentication (i.e., username and password). It is not possible to infer these values by either static or dynamic analysis, or by randomly guessing. Therefore, our tool was designed to use a prespecified set of values for database authentication, e.g., User = joe, Password = "12345." Currently, the only values we need to specify are username/password pairs for user authentication. From this information, Apollo has been able to extract all other required information from the database.

## 3.1 Algorithm

Fig. 2 shows pseudocode for our algorithm. The inputs to the algorithm are: a program $\mathcal{P}$, an oracle for the output $\mathcal{O}$, and an initial state of the environment $\mathcal{S}_0$. The output of the algorithm is a set of bug reports $\mathcal{B}$ for the program $\mathcal{P}$, according to $\mathcal{O}$. Each report consists of a single failure, defined by the error message and the set of statements that is related to the failure. In addition, the report contains the set of all inputs under which the failure was exposed, and the set of all path constraints that lead to the inputs exposing the failure.

The algorithm uses a queue of configurations. Each configuration is a pair of a path constraint and an input. A *path constraint* is a conjunction of conditions on the program's input parameters. The queue is initialized with the empty path constraint and the empty input (line 3). The program is executed concretely on the input (line 6) and tested for failures by the oracle (line 7). Then, the path constraint and input for each detected failure are merged into the corresponding bug report (lines 7 and 8).

Next, the algorithm uses a subroutine, *getConfigs*, to find new configurations. First, the program is executed symbolically on the same input (line 15). The result of symbolic

**parameters**: Program $\mathcal{P}$, oracle $O$, Initial state $\mathcal{S}_0$
**result**       : Bug reports $\mathcal{B}$;
$\mathcal{B}$ : $setOf(\langle$failure, $setOf(\text{pathConstraint}), setOf(\text{input})\rangle)$

1   $\mathcal{B} := \varnothing$;
2   $toExplore := emptyQueue()$;
3   $enqueue(toExplore, \langle emptyPathConstraint(), emptyInput\rangle)$;
4   **while** $not\ empty(toExplore)\ and\ not\ timeExpired()$ **do**
5      $\langle pathConstraint, input\rangle := dequeue(toExplore)$;
6      $output := executeConcrete(\mathcal{S}_0, \mathcal{P}, input)$;
7      **foreach** $f$ $in$ $getFailures(O, output)$ **do**
8         merge $\langle f, pathConstraint, input\rangle$ into $\mathcal{B}$;
9      $newConfigs := getConfigs(input)$;
10     **foreach** $\langle pathConstraint_i, input_i\rangle \in newConfigs$ **do**
11        $enqueue(toExplore, \langle pathConstraint_i, input_i\rangle)$;
12   **return** $\mathcal{B}$;

13   **Subroutine** $getConfigs(input)$:
14   $configs := \varnothing$;
15   $c_1 \wedge \ldots \wedge c_n := executeSymbolic(\mathcal{S}_0, \mathcal{P}, input)$;
16   **foreach** $i = 1, \ldots, n$ **do**
17      $newPC := c_1 \wedge \ldots \wedge c_{i-1} \wedge \neg c_i$;
18      $input := solve(newPC)$;
19      **if** $input \neq \perp$ **then**
20         $enqueue(configs, \langle newPC, input\rangle)$;
21   **return** $configs$;

Fig. 2. The failure detection algorithm. The output of the algorithm is a set of bug reports. Each bug report contains a failure, a set of path constraints exposing the failure, and a set of inputs exposing the failure. The $solve$ auxiliary function uses the constraint solver to find an input satisfying the path constraint or returns $\perp$ if no satisfying input exists. The $merge$ auxiliary function merges the pair of pathConstraint and input for an already detected failure into the bug report for that failure.

**parameters**: Program $\mathcal{P}$, oracle $O$, bug report $b$
**result**       : Short path constraint that exposes $b.failure$

1   $c_1 \wedge \ldots \wedge c_n := intersect(b.pathConstraints)$;
2   $pc := true$;
3   **foreach** $i = 1, \ldots, n$ **do**
4      $pc_i := c_1 \wedge \ldots c_{i-1} \wedge c_{i+1} \wedge \ldots c_n$;
5      **if** $!exposesFailures(pc_i)$ **then**
6         $pc := pc \wedge c_i$;
7   **if** $exposesFailures(pc)$ **then**
8      **return** $pc$;
9   **return** $shortest(b.pathConstraints)$;

10   **Subroutine** $exposesFailure(pc)$:
11   $input_{pc} := solve(pc)$;
12   **if** $input_{pc} \neq \perp$ **then**
13      $output_{pc} := executeConcrete(\mathcal{P}, input_{pc})$;
14      $failures_{pc} := getFailures(O, output_{pc})$;
15      **return** $b.failure \in failures_{pc}$;
16   **return** $false$;

Fig. 3. The path constraint minimization algorithm. The method $intersect$ returns the set of conjuncts that are present in all given path constraints, and the method $shortest$ returns the path constraint with fewest conjuncts. The other auxiliary functions are the same as in Fig. 2.

The algorithm now enters the **foreach** loop on line 16 of Fig. 2, and starts generating new path conditions by systematically traversing subsequences of the above path constraint, and negating the last conjunct. Hence, from (I), the algorithm derives the following three path constraints:

$$NotSet(\text{page}) \wedge \text{page2} \neq 1337 \wedge \text{login} = 1 \qquad \text{(II)}$$
$$NotSet(\text{page}) \wedge \text{page2} = 1337 \qquad \text{(III)}$$
$$Set(\text{page}) \qquad \text{(IV)}$$

**Iteration 2.** For path constraint (II), the constraint solver may find the following input (the solver is free to select any value for `page2` other than 1337): `page2` $\leftarrow \emptyset$, `login` $\leftarrow 1$.

When the program is executed with this input, the condition of the if-statement on line 16 evaluates to `true`, resulting in a call to the `validateLogin` method. Then, the condition of the if-statement on line 28 evaluates to `true` because the `username` parameter is not set, resulting in the generation of output containing an incorrect HTML tag `j2` on line 29. When the HTML validator checks the page, the failure is discovered and a bug report is created and added to the output set of bug reports.

### 3.3 Path Constraint Minimization

The failure detection algorithm (Fig. 2) returns bug reports. Each bug report contains a set of path constraints, and a set of inputs exposing the failure. Previous dynamic test generation tools [7], [17], [39] presented the whole input (i.e., many $\langle inputParameter, value\rangle$ pairs) to the user without an indication of the subset of the input responsible for the failure. As a postmortem phase, our minimization algorithm attempts to find a shorter path constraint for a given bug report (Fig. 3). This eliminates irrelevant constraints, and a solution for a shorter path constraint is often a smaller input.

For a given bug report $b$, the algorithm first intersects all the path constraints exposing $b.failure$ (line 1). The minimizer systematically removes one conjunct at a time (lines 3-6). If one of these shorter path constraints does not

execution is a path constraint, $\bigwedge_{i=1}^{n} c_i$, that is satisfied by the path that was just executed from entry to exit of the whole program. The subroutine then creates new inputs by solving modified versions of the path constraint (lines 16-20), as follows: For each prefix of the path constraint, the algorithm negates the last conjunct (line 17). A solution, if it exists, to such an alternative path constraint corresponds to an input that will execute the program along a prefix of the original execution path, and then take the opposite branch, presumably covering new code. The algorithm uses a constraint solver to find a concrete input for each path constraint (line 18).

### 3.2 Example

Let us now consider how the algorithm of Fig. 2 exposes the third fault in the example program of Fig. 1.

**Iteration 1.** The first input to the program is the empty input, which is the result of solving the empty path constraint. During the execution of the program on the empty input, the condition on line 6 evaluates to `true`, and `page` is set to $\emptyset$. The condition on line 10 evaluates to `false`. The condition on line 16 evaluates to `false` because parameter `login` is not defined. The `switch` statement on line 18 selects the case on line 20 because `page` has the value of 0. Execution terminates on line 26. The HTML verifier determines that the output is legal, and $executeSymbolic$ produces the following path constraint:

$$NotSet(\text{page}) \wedge \text{page2} \neq 1337 \wedge \text{login} \neq 1 \quad \text{(I)}$$

expose $b.failure$, then the removed conjunct is required for exposing $b.failure$. The set of all such required conjuncts determines the minimized path constraint. From the minimized path constraint, the algorithm produces a concrete input that exposes the failure.

The algorithm in Fig. 3 does not guarantee that the returned path constraint is the shortest possible that exposes the failure. However, the algorithm is simple, fast, and effective in practice (see Section 6.3.2).

Our approach for minimization is similar in spirit to delta debugging [10], [47], a well-known input minimization technique. However, our algorithm operates on the *path constraint* that exposes the failure, and not directly on the *input*. A constraint concisely describes a class of inputs (e.g., the constraint page2 $\neq$ 1337 describes all inputs different than 1337). Since a concrete input is an instantiation of a constraint, it is more effective to reason about input properties in terms of their constraints. Our current minimization algorithm also differs from delta debugging in that it does not rely on efficient binary search like techniques to identify redundant components of path constraints (thus far, the identified path constraints were quite small in practice, and the use of such techniques was unnecessary).

Each failure might be encountered along several execution paths that might partially overlap. Without any information about the properties of the inputs, delta debugging minimizes only a *single* input at a time, while our algorithm handles *multiple* path constraints that lead to a failure.

### 3.4 Minimization Example

The malformed HTML failure described in Section 3.2 can be triggered along different execution paths. For example, both of the following path constraints lead to inputs that expose the failure. Path constraint (*a*) is the same as (II) in Section 3.2.

$$NotSet(\text{page}) \land \text{page2} \neq 1337 \land \text{login} = 1 \qquad (a)$$
$$Set(\text{page}) \land \text{page} = \emptyset \land \text{page2} \neq 1337 \land \text{login} = 1 \quad (b)$$

First, the minimizer computes the intersection of the path constraints (line 1). The intersection is:

$$\text{page2} \neq 1337 \land \text{login} = 1 \quad (a \cap b)$$

Then, the minimizer creates two shorter path constraints by removing each of the two conjuncts in turn. First, the minimizer creates path constraint login = 1. This path constraint corresponds to an input that reproduces the failure, namely login ← 1. The minimizer determines this by executing the program on the input (line 14 in Fig. 3). Second, the minimizer creates path constraint page2 $\neq$ 1337. This path constraint does not correspond to an input that exposes the failure. Thus, the minimizer concludes that the condition login = 1, that was removed from $(a \cap b)$ to form the second path constraint, is required. In this example, the minimizer returns login = 1. The result is the minimal path constraint that describes the minimal failure-inducing input, namely, login ← 1.

## 4 COMBINED CONCRETE AND SYMBOLIC EXECUTION WITH EXPLICIT-STATE MODEL CHECKING

A typical PHP Web application is a client-server application in which data and control flows interactively between a server that runs PHP scripts and a client, which is usually a Web browser. The PHP scripts that run on the server generate HTML that includes interactive user-input widgets such as buttons and menu items that, when selected by the user, invoke other PHP scripts. When these other PHP scripts are invoked, they are passed a combination of user input and constant values taken from the generated HTML. Modeling such user input is important because coverage of the application will typically remain very low otherwise.

In Section 3, we described how to find failures in PHP Web applications by adapting an existing test generation approach to consider language constructs that are specific to PHP, by using an oracle to validate the output, and by having the user supply values that are needed for database authentication. However, in practice, a PHP application starts by executing a script, which generates an HTML page from which other PHP scripts may be invoked as a result of user actions (e.g., by selecting a menu entry). Such additional scripts may refer to: 1) parameters that are transferred as part of the call, 2) session state that is persisted in the environment, and 3) to the database (after it has been updated as a result of executing a previous script). The solution that we previously presented does not handle these issues, and is therefore incapable of achieving good coverage for realistic PHP applications.

To handle this problem, Apollo implements a form of explicit-state software model checking. That is, Apollo systematically explores the state space of the system, i.e., the program under test. The algorithm in Section 3 always restarts the execution from the same initial state and discards the state reached at the end of each execution. Thus, the algorithm reaches only one-level deep into the application, where each level corresponds to a cycle of: a PHP script that generates an HTML form that the user interacts with to invoke the next PHP script. In contrast, the algorithm presented in this section remembers and restores the state between executions of PHP scripts. This technique, known as state matching, is widely known in model checking [22], [42] and implemented in tools such as SPIN [13] and JavaPath-Finder [21]. To our knowledge, we are the first to implement state matching in the context of Web applications and PHP.

### 4.1 Interactive User Simulation Example

Fig. 4 shows an example of a PHP application that is designed to illustrate the particular complexities of finding faults in an interactive web application. In particular, the figure shows: an index.php top-level script that contains static HTML in Fig. 4a, a generic login script login.php in Fig. 4c, and a skeleton of a data display script view.php in Fig. 4d. The PHP scripts in Fig. 4 rely on a shared include file constants.php that defines some standard constants, which is shown in Fig. 4b. Note that the code in Fig. 4 is an ad-hoc mixture of PHP statements and HTML fragments. The PHP code is delimited by <?php and ?> tokens (see, e.g., lines 44 and 69 in Fig. 4d). The use of HTML in the middle of PHP code indicates that HTML is generated as if it were the argument of a print statement. The dirname function—which returns the directory component of a filename—is used in the require statements as an example of including a file whose name is computed at runtime.

These PHP scripts are part of the client-server work flow in a Web application: The user first sees the index.php

```
1   <html>                                          10   <?php
2   <head>Login</head>                              11     userTag = 'user'
3   <body>                                           12     pwTag = 'pw';
4     <form name="login" action="login.php">        13     typeTag = 'type';
5       <input type="text" name="user"/>            14   ?>
6       <input type="password" name="pw"/>
7     </form>
8   </body>
9   </html>                  (a)                                          (b)


15  <HTML>                                           42   <HTML>
16  <?php                                            43   <HEAD>Topic View</HEAD>
17    require( dirname(__FILENAME__).'/includes/constants.php');   44   <?php
                                                     45     print "<BODY>\n";
18                                                   46     if(check_password($_SESSION[$userTag], $_SESSION[$pwTag]) {
19    $user = $_REQUEST[ 'user' ];                   47       require( dirname(__FILENAME__).'/includes/constants.php');
20    $pw = $_REQUEST[ 'pw' ];                       48
21                                                   49       $type = $_SESSION[ $typeTag ];
22    if (check_password($user, $pw) {               50       $topic = $_REQUEST[ 'topic' ];
23      print "<HEAD>Login Successful</HEAD>\n";     51
24      $_SESSION[ $userTag] = $user;                52       if ($type == 'admin') {
25      $_SESSION[ $pwTag ] = $pw;                   53         print "<H1>Admin ";
26    }                                              54       } else {
27    else {                                         55         print "<H1>Normal ";
28      print "<HEAD>Login Failed</HEAD>\n";         56       }
29    }                                              57       print "View of $topic</H1>\n";
30  ?>                                               58
31    <BODY>                                         59       /* code to print topic view... */
32      <FORM action="view.php">                     60
33        <INPUT TYPE="text" NAME="topic"/>          61       if ($type == 'admin') {
34      </FORM>                                       62         print "<H2>Administrative Details\n";
35    </BODY>                                         63         /* code to print admin details... */
36  <?php                                            64       }
37    if ($user == 'admin') {                        65     } else {
38      $_SESSION[ $typeTag ] = 'admin';             66       print "Please Log in\n";
39    }                                              67     }
40  ?>                                               68     print "</BODY>\n";
41  </HTML>                                          69   ?>
                                                     70   </HTML>
                         (c)                                              (d)
```

Fig. 4. Example of a PHP Web application. (a) `index.php`, (b) `constants.php`, (c) `login.php`, and (d) `view.php`.

page of Fig. 4a and enters credentials. The user-input credentials are processed by the script in Fig. 4c, which generates a response page that allows the user to enter further input—a topic—that in turn entails further processing by the script in Fig. 4d. Note that the username and password that are entered by the user during the execution of `login.php` are stored in special locations $_SESSION [ $userTag ] and $_SESSION[ $pwTag ], respectively. Moreover, if the user is the administrator, this fact is recorded similarly, in $_SESSION[ $typeTag ]. These locations illustrate how PHP handles *session state*, which are data that persist from one page to another, typically for a particular interaction by a particular user. Thus, the updates to _SESSION in Fig. 4c will be seen (as the SESSION information is saved and read locally on the server) by the code in Fig. 4d when the user follows the link to `view.php` in the HTML page that is returned by `login.php`. The `view.php` script uses this session information to verify the username/password in line 46.

Our example program contains an error in the HTML produced for the administrative details: The H2 tag that is opened on line 62 of Fig. 4d is not closed. While this fault itself is trivial, finding it is not. Assume that testing starts (as an ordinary user would) by entering credentials to the script in Fig. 4c. A tester must then discover that setting $user to the value "admin" results in the selection of a different branch that records the user type "admin" in the session state (see lines 37-39 in `login.php`). After that, a tester would have to enter a topic in the form generated by the login script, and would then proceed to Fig. 4d with the appropriate session state, which will finally generate HTML

exhibiting the fault as is shown in Fig. 5a. Thus, finding the fault requires a careful selection of inputs to a series of interactive scripts, as well as making sure that updates to the session state during the execution of these scripts are preserved (i.e., making sure that the executions of the different scripts happen during the same session).

## 4.2 Algorithm

Fig. 6 shows pseudocode for the algorithm, which extends the algorithm in Fig. 2 with explicit-state model checking to handle the complexity of simulating user inputs. The algorithm tracks the state of the environment, and automatically discovers additional configurations based on an analysis of the output for available user options. In particular, the algorithm 1) tracks changes to the state of the environment (i.e., session state, cookies, and the database) and 2) performs an "on-the-fly" analysis of the output produced by the program to determine what user options it contains, with their associated PHP scripts. By determining the state of the environment as it exists when an HTML page is produced, the algorithm can determine the environment in which additional scripts are executed as a result of user interaction. This is important because a script is much more likely to perform complex behavior when executed in the correct context (environment). For example, if the Web application does not record in the environment that a user is logged in, most subsequent calls will terminate quickly (e.g., when the condition in line 46 of Fig. 4d is false) and will not present useful information. For simplicity, the algorithm implicitly handles the fact that there are possibly multiple entry points

```
1 <HTML>
2 <HEAD>Topic View</HEAD>
3 <BODY>
4 <H1>Admin View of A topic</H1>
 ...
5 <H2>Administrative Details
 ...
6 </BODY>        (a)
7 </HTML>

 Error at line 6, character 7:  end tag for "H2" omitted; possible causes include a missing
 end tag, improper nesting of elements, or use of an element where it is not allowed
 Line 5, character 1:  start tag was here

                                          (c)
```

| HTML line | PHP lines in 4(d) |
|---|---|
| 1 | 42 |
| 2 | 43 |
| 3 | 45 |
| 4 | 53, 57 |
| 5 | 62 |
| 6 | 68 |
| 7 | 70        (b) |

Fig. 5. (a) HTML produced by the script of Fig. 4d. (b) Output mapping constructed during execution. (c) Part of output of WDG Validator on the HTML of (a).

into a PHP program. Thus, an input will contain the script to execute in addition to the values of the parameters. For instance, the first call might be to the index.php script, while subsequent calls can execute other scripts.

There are four differences (underlined in the figure) with the simplified algorithm that was previously shown in Fig. 2.

1. A configuration contains an explicit state of the environment (before the only state that was used was the initial state $\mathcal{S}_0$) in addition to the path constraint and the input (line 3).

2. Before the program is executed, the algorithm (method executeConcrete) will restore the environment to the state given in the configuration (line 7) and will return the new state of the environment after the execution.

3. When the getConfigs subroutine is executed to find new configurations, it analyzes the output to find possible transitions from the new environment state (lines 24-27). The analyzeOutput function extracts parameter names and possible values for each parameter, and represents the extracted information as a path constraint. For simplicity, the algorithm uses only one entry point into the program. However, in practice, there maybe several entry points into the program (e.g., it is possible to call different PHP scripts). The analyzeOutput function discovers these entry points in addition to the path constraints. In practice, each transition is expressed as a pair of a path constraint and an entry point.

4. The algorithm uses a set of configurations that are already in the queue (line 14) and it performs state matching in order to only explore new configurations (line 11).

### 4.3 Example

We will now illustrate the algorithm of Fig. 6 using the example application of Fig. 4. The inputs to the algorithm are: $\mathcal{P}$ is the code from Fig. 4, the initial state of the environment is empty, the first script to execute is the script in Fig. 4a, and $\mathcal{O}$ is the WDG HTML validator.[8] The algorithm begins on line 3 by initializing the work queue with one item: an empty input to the script of Fig. 4a with an empty path constraint and an empty initial environment.

8. http://htmlhelp.com/tools/validator/.

**Iteration 1.** The first iteration of the outer loop (lines 5-14) removes that item from the queue (line 6), restores the empty initial state, and executes the script (line 7).

```
parameters: Program 𝒫, oracle O, Initial state 𝒮₀
result      : Bug reports ℬ;
ℬ : setOf(⟨failure, setOf(pathConstraint), setOf(input)⟩)
1  ℬ := ∅;
2  toExplore := emptyQueue();
3  enqueue(toExplore, ⟨emptyPC(), emptyInput(), 𝒮₀⟩);
4  visited := {⟨emptyPathConstraint(), emptyInput(), 𝒮₀⟩};
5  while not empty(toExplore) and not timeExpired() do
6     ⟨pathConstraint, input, 𝒮_start⟩ := dequeue(toExplore);
7     ⟨output, 𝒮_end⟩ := executeConcrete(𝒮_start, 𝒫, input);
8     foreach f in getFailures(O, output) do
9        merge ⟨f, pathConstraint, input⟩ into ℬ;
10    newConfigs := getConfigs(input, output, 𝒮_start, 𝒮_end);
11    newConfigs := newConfigs − visited;
12    foreach ⟨pathConstraint_i, input_i, 𝒮_i⟩ ∈ newConfigs do
13       enqueue(toExplore, ⟨pathConstraint_i, input_i, 𝒮_i⟩);
14       visited := visited ∪ {⟨pathConstraint_i, input_i, 𝒮_i⟩};
15 return ℬ;

16 Subroutine getConfigs(input, output, 𝒮_start, 𝒮_end):
17 configs := ∅;
18 c₁ ∧ ... ∧ cₙ := executeSymbolic(𝒮_start, 𝒫, input);
19 foreach i = 1,...,n do
20    newPC := c₁ ∧ ... ∧ c_{i−1} ∧ ¬c_i;
21    input := solve(pathConstraint);
22    if input ≠ ⊥ then
23       enqueue(configs, ⟨newPC, input, 𝒮_start⟩);
24    foreach newPC_i ∈ analyzeOutput(output) do
25       newInput := solve(newPC_i);
26       if newInput ≠ ⊥ then
27          configs := configs ∪ ⟨newPC_i, newInput_i, 𝒮_end⟩;
28 return configs;
```

Fig. 6. The failure detection algorithm. The output of the algorithm is a set of bug reports; each reports a failure and the set of tests exposing that failure. The *solve* auxiliary function uses the constraint solver to find an input satisfying the path constraint, or returns ⊥ if no satisfying input exists. The *merge* auxiliary function merges the pair of pathConstraint and input for an already detected failure into the bug report for that failure. The *analyzeOutput* auxiliary function performs an analysis of the output to extract possible transitions from the current environment state.
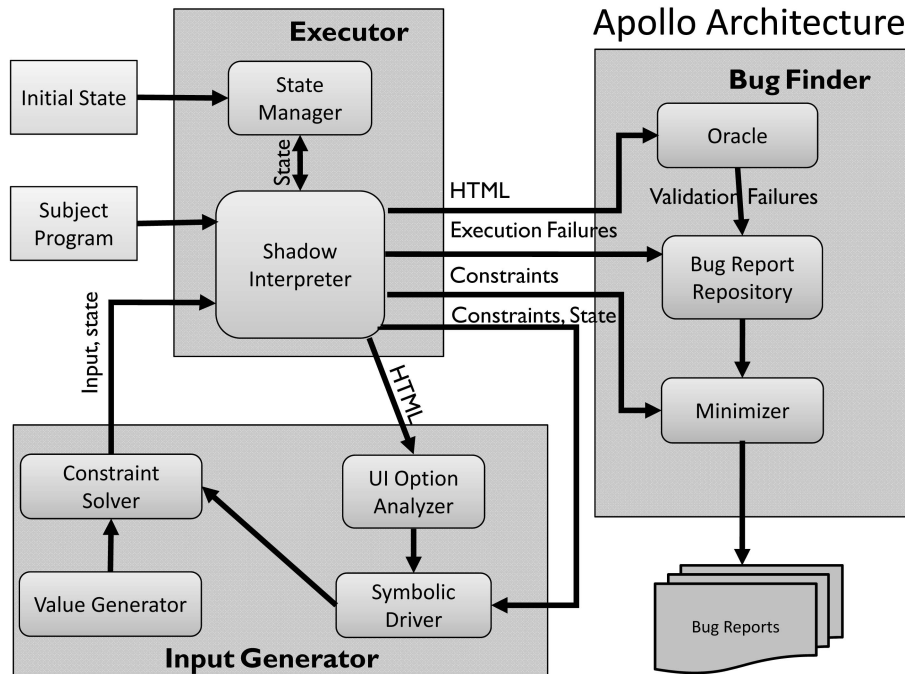
Fig. 7. The architecture of Apollo.

No failures are observed. The call to *executeSymbolic* on line 18 returns an empty path constraint, so the function *analyzeOutput* on line 24 is executed next, and returns one user option; $\langle$login.php, $\emptyset, \emptyset\rangle$ for executing login.php with no input, and the empty state. This configuration is added to the queue (line 13) since it was not seen before.

**Iteration 2-5.** The next iteration of the top-level loop dequeues the new work item, and executes login.php with empty input, and empty state. No failures are found. The call to *executeSymbolic* in line 18 returns a path constraint user $\neq$ admin $\wedge$ user $\neq$ reg, indicating that the call to check_password on line 22 in Fig. 4c returned false.[9] Given this, the loop at lines 19-23 will generate several new work items for the same script with the following path constraints: user $\neq$ admin $\wedge$ user $=$ reg and user $=$ admin, which are obtained by negating the previous path constraint. The loop on lines 24-27 is not entered because no user-input options are found. After several similar iterations, two inputs are discovered: user $=$ admin $\wedge$ pw $=$ admin and user $=$ reg $\wedge$ pw $=$ reg. These corresponds to alternate control flows in which the check_password test succeeds.

**Iteration 6-7.** The next iteration of the top-level loop dequeues an item. Using this item, the call to check_password will succeed (assume it selected user $=$ reg...). Once again, no failures are observed, but now the session state with *user* and *pw* set is recorded at line 7. Also, this time *analyzeOutput* (line 24) finds the link to the script in Fig. 4d, and so the loop at lines 24-27 adds one item to the queue, executing view.php with the current session state.

The next iteration of the top-level loop dequeues one work item. Assume that it takes the last one described above. Thus, it executes the script in Fig. 4d with a session

that defines *user* and *pw* but not *type*. Hence, it produces an execution with no errors.

**Iteration 8-9.** The next loop iteration takes that last work item, containing a user and password pair for which the call to check_password succeeds, with the username as "admin." Once again, no failures occur, but now the session state with *user*, *pw*, and *type* set is recorded at line 7. This time, there are no new inputs to be derived from the path constraint since all prefixes have been covered already. Once again, parsing the output finds the link to the script in Fig. 4d and adds a work item to the queue, but with a different session state (in this case, the session state also includes a value for *type*). The resulting execution of the script in Fig. 4d with the session state that includes *type* results in an HTML failure.

## 5 IMPLEMENTATION

We created a tool called Apollo that implements our technique for PHP. Apollo consists of three major components, **Executor**, **Bug Finder**, and **Input Generator** illustrated in Fig. 7. This section first provides a high-level overview of the components and then discusses the pragmatics of the implementation.

The inputs to Apollo are the program under test and an initial value for the environment. The environment of a PHP program consists of the database, cookies, and stored session information. The initial environment usually consists of a database populated with some values, and user-supplied information about username/password pairs to be used for database authentication.[10]

---

9. For simplicity, we omit the details of this function. It compares the username and password to some constants "admin" and "reg."

10. Attempting to retrieve information from the database using randomly chosen values for username/password is unlikely to be successful. Symbolic execution is equally helpless without the database manager because reversing cryptographic functions is beyond the state of the art for constraint solvers.

The **Executor** is responsible for executing a PHP script with a given input in a given state. The executor contains two subcomponents:

- The **Shadow Interpreter** is a PHP interpreter that we have modified to propagate and record path constraints and positional information associated with output. This positional information is used to determine which failures are likely to be symptoms of the same fault.
- The **State Manager** restores the given state of the environment (database, session, and cookies) before the execution and stores the new environment after the execution.

The **Bug Finder** uses an oracle to find HTML failures, stores all bug reports, and finds the minimal conditions on the input parameters for each bug report. The Bug Finder has the following subcomponents:

- The **Oracle** finds HTML failures in the output of the program.
- The **Bug Report Repository** stores all bug reports found during all executions.
- The **Input Minimizer** finds, for a given bug report, the smallest path constraint on the input parameters that results in inputs inducing the same failure as in the report.

The **Input Generator** implements the algorithm described in Fig. 6. The Input Generator contains the following subcomponents:

- The **UI Option Analyzer** analyzes the HTML output of each execution to convert the interactive user options into new inputs to execute.
- The **Symbolic Driver** generates new path constraints from the constraints found during the execution.
- The **Constraint Solver** computes an assignment of values to input parameters that satisfies a given path constraint.
- The **Value Generator** generates values for parameters that are not otherwise constrained, using a combination of random value generation and constant values mined from the program source code.

## 5.1 Executor

We modified the Zend PHP interpreter 5.2.2[11] to produce symbolic path constraints for the executed program, using the "shadow interpreter" approach [11]. The shadow interpreter performs the regular (concrete) program execution using the concrete values, and simultaneously performs symbolic execution. Creating the shadow interpreter required five alterations to the PHP runtime:

1. **Associating Symbolic Parameters with Values**. Conceptually, we associate a symbolic variable with every value that arises at runtime; however, we are interested in tracking uses of input values, and so the only variables that contain nonempty values are those for the inputs themselves and values computed from them in one of the following ways: variable assignment, parameter passing, string concatenation,

and conversion to a number. Clearly, this potentially leaves many uses of input unaccounted for. However, our results suggest that this is sufficient to capture the bulk of how PHP code uses inputs in practice. Values derived directly from input are those read from one of the special arrays _POST, _GET, and _REQUEST, which store parameters supplied to the PHP program. For example, executing the statement $x = \$\_GET["param1"]$ results in associating the value read from the global parameter param1 and bound to parameter x with the symbolic variable param1. Values maintain their associations through the operations mentioned above; that is, the symbolic variables for the new values receive the same value as the source value had. Importantly, during program execution, the concrete values remain, and the shadow interpreter does not influence execution.

Unlike other projects that perform concrete and symbolic execution [7], [17], [18], [39], our interpreter does not associate complex symbolic expressions with all runtime values, but only symbolic variables, which exist only for input-derived values. This design keeps the constraint solver simple and reduces the performance overhead. As our results (Section 6) indicate, this lightweight approach is sufficient for the analyzed PHP programs.

2. **Storing Constraints at Branch Points**. At branching points (i.e., value comparisons) that involve values associated with symbolic variables, the interpreter extends the initially empty path constraint with a conjunct that corresponds to the branch actually taken in the execution. For example, if the program executes a statement if ($name == "John") and this condition succeeds, where $name is associated with the symbolic variable username, then the algorithm appends the conjunct username = "John" to the path constraint.

3. **Handling PHP Native Functions**. Our modified interpreter records conditions for PHP-specific comparison operations, such as isset and empty, which can be applied to any variable. Operation isset returns a boolean value that indicates whether or not a value different from NULL was supplied for a variable. The empty operator returns true when applied to: the empty string, $\emptyset$, "$\emptyset$," NULL, false, or an empty array. The interpreter records the use of isset on values with an associated symbolic variable, and on uninitialized parameters.

The isset comparison creates either the *NotSet* or the *Set* condition. The constraint solver chooses an arbitrary value for a parameter p if the only condition for p is *Set* (p). Otherwise, it will also take into account other conditions. The *NotSet* condition is used only in checking the feasibility of a path constraint. A path constraint with the *NotSet* (p) condition is feasible only if it does not contain any other conditions on p. The empty comparison creates equality or inequality conditions between the parameter and the values that are considered empty by PHP.

4. **Propagating Inputs through Sessions and Cookies**. While HTTP is a stateless protocol, various mechanisms are used to thread a series of HTTP requests

into a transaction, a major one being server-side session state. In PHP, this is exposed directly as the _SESSION variable, and the interpreter propagates this persistent state across multiple requests that are deemed part of the same session by the server.

The use of session state allows a PHP application to store user-supplied information on the server for retrieval by other scripts. We enhanced the PHP interpreter to record when input parameters are stored in session state. This enables Apollo to track constraints on input parameters in all scripts that use them.

5. **Web Server Integration**. Dynamic Web applications often depend on information supplied by a Web server (such as Apache), and some PHP constructs are simply ignored by the command line interpreter (e.g., *header*). In order to allow Apollo to analyze more PHP code, Apollo supports execution through the Apache Web server in addition to the stand-alone command line executor. A developer can use Apollo to silently analyze the execution and record any failure found while manually using the subject program on an Apache server.

The modified interpreter performs symbolic execution along with concrete execution, i.e., every variable during program execution has a concrete value and may have additionally a symbolic value. Only the concrete values influence the control flow during the program execution, while the symbolic execution is only a "witness" that records, but does not influence, control-flow decisions at branching points. This design deals with exceptions naturally because exceptions do not disrupt the symbolic-value mapping for variables.

Our approach to symbolic execution allows us to handle many PHP constructs that are problematic in a purely static approach. For instance, for computed variable names (e.g., $x =${$foo}), any symbolic information associated with the value that is held by the variable named by foo will be passed to x by the assignment.[12] In order to heuristically group HTML failures that may be manifestations of the same fault, Apollo records the output statement (i.e., echo or print) that generated each fragment of HTML output.

Retargeting the existing PHP interpreter for concolic execution required the following changes to the interpreter:

1. Extending the values manipulated during execution (zval in the case of the Zend implementation of the PHP interpreter) to track symbolic expressions.
2. Instrumenting user-input access to start the symbolic tracking or to note that the program queries about a missing input parameter. In the case of the Zend PHP interpreter, this was done by instrumenting accesses to the global parameter tables.
3. Detecting interesting execution states, such exceptions, and unclean exits.
4. Instrumenting all internal calls responsible for output, in order to track the relations between output and executed statements.

5. Instrumenting the calls to read and write sessions and cookies, and changing the representation of the stored value in order to store symbolic expressions in the environment.

**State manager.** PHP applications make use of persistent state such as the database, session information, and cookies. The State Manager is in charge of 1) restoring the environment prior to each execution and 2) storing the new environment after each execution.

## 5.2 Bug Finder

The bug finder is in charge of transforming the results of the executed inputs into bug reports. Below is a detailed description of the components of the bug finder.

**Bug report repository.** This repository stores the bug reports found in all executions. Each time a failure is detected, the corresponding bug report (if the same failure was discovered before) is updated with the path constraint and the configuration inducing the failure. A failure is uniquely defined by the following set of characteristics: the type of the failure (execution failure or HTML failure), the corresponding message (PHP error/warning message for execution failures and validator message for HTML failures), and the PHP statement generating the problematic HTML fragments identified by the validator (for HTML failures) or the PHP statement involved in the PHP interpreter error report (for execution failures). When the exploration is complete, each bug report contains one failure characteristic (error message and statement involved in the failure) and the sets of path constraints and inputs exposing failures with the same characteristics. Recording the constraints under which a failure occurs can be helpful for debugging because the state could be large or complex, and an error might only be reproducible in a certain state.

**Oracle.** PHP Web applications output HTML/XHTML. Therefore, in Apollo, we use as oracle an HTML validator that returns syntactic (malformed) HTML failures found in a given document. We experimented with both the offline WDG validator[13] and the online W3C markup validation service.[14] Both oracles identified the same HTML failures. Our experiments use the faster WDG validator.

**Input minimizer.** Apollo implements the algorithm described in Fig. 3 to perform *postmortem* minimization of the path constraints. For each bug report, the minimizer executes the program multiple times, with multiple inputs that satisfy different path constraints, and attempts to find the shortest path constraint that results in the same failure characteristics.

## 5.3 Input Generator

**UI option analyzer**. Many PHP Web applications create interactive HTML pages that contain user-interface elements such as buttons and menus that allow the user interaction needed to execute further parts of the application. In such cases, pressing the button may result in the execution of additional PHP source files. There are two challenges involved in dealing with such interactive applications.

---

12. On the other hand, any data flow that passes outside PHP, such as via JavaScript code in the generated HTML, will not be tracked by this approach.

13. http://htmlhelp.com/tools/validator/offline.
14. http://validator.w3.org.

```php
<?php
   echo "<h2>WebChess ".$Version." Login"</h2>;
?>
<form method="post" action="mainmenu.php">
<p>
   Nick: <input name="txtNick" type="text" size="15" default="admin"/>
   <br />
   Password: <input name="pwdPassword" type="password" size="15"/>
</p>
<p>
   <input name="login" value="login" type="submit"/>
   <input name="newAccount" value="New Account"
     type="button" onClick="window.open('newuser.php', '_self')"/>
</p>
</form>
```

Fig. 8. A simplified version of the main entry point (`index.php`) to a PHP program. The HTML output of this program contains a form with two buttons. Pressing the `login` button executes `mainmenu.php` and pressing the `newAccount` button will execute the `newuser.php` script.

- The HTML output must be examined to find the referenced scripts and the different values that can be supplied as parameters.
- Apollo needs to be able to follow input parameters through the shared global information (database, the session, and the cookie mechanisms).

Apollo's approach to the above challenges is to simulate user interaction by analyzing the dynamically created HTML output and tracking the symbolic parameters through the environment:

- Apollo automatically extracts the available user options from the HTML output so that it collects all HTML forms in the page and their components, e.g., buttons and text areas, through which the user can provide input. Any default values for such elements are also collected.
- Apollo collects static HTML documents that can be called from the dynamic HTML output, i.e., Apollo gather all href attributes in the HTML document.
- Apollo performs a cursory analysis of JavaScript code to find other syntactic references, for instance, a window.open call with a static url as a parameter.

Since additional code on the client side (for instance, JavaScript) might be executed when a button is pressed, this approach might induce false positive bug reports. In our experiments, this limitation produced no false positive bug reports.

For example, after analyzing the output of the program of Fig. 8, the UI Option Analyzer will return the following two options:

1. Script: "mainmenu.php"
    PathConstraint:

$$\texttt{txtNick} = {}^{\prime\prime}Admin^{\prime\prime} \wedge \text{Exist} \ (\texttt{pwdPassword})$$

2. Script: "newuser.php"
    PathConstraint: $\emptyset$

For instance, the first of the items in the above list comes from the submit option to the form in Fig. 8: The mainmenu.php script is directly specified by the submit option, and the path constraint comes from the input items in the form. Any string may be provided for pwdPassword, so we generate a constraint just that it exist. But the txtNick item is given the default value since that input is provided

directly by the form. Note that further exploration can cause either of these constraints to be negated, which will simulate providing other possible inputs. The second of these items comes from onClick = "window.open('newuser.php', '_self'"', which directly opens a new window pointing to the newuser.php script.

The **Symbolic Driver** implements the combined concrete and symbolic algorithm of Fig. 2. The driver has two main tasks: Select which input to consider next (line 5), and create additional inputs from each executed input (by negating conjuncts in the path constraint). To select which input to consider next, the driver uses a *coverage heuristic*, similar to those used in EXE [7] and SAGE [18]. Each conjunct in the path constraint knows the branch that created the conjunct, and the driver keeps track of all branches previously executed and favors inputs created from path constraints that contain unexecuted branches.

To avoid redundant exploration of similar executions, Apollo performs state matching (performed implicitly in Line 11 of Fig. 6) by not adding already explored transitions.

**Constraint solver.** The interpreter implements a lightweight symbolic execution, in which the only constraints are equality and inequality with constants. Apollo transforms path constraints into integer constraints in a straightforward way, and uses choco[15] to solve them.

This approach still allows us to handle values of the standard types (integer, string), and is straightforward because the only constraints are equality and inequality.[16]

In cases where parameters are unconstrained, Apollo randomly chose values from a predefined list of constants.

While limiting to the basic types number and string and only comparisons may seem very restrictive, note that all input comes to PHP as strings; furthermore, in our experience, the bulk of use of input values consists of the kinds of simple operations that are captured by our tracing and the kinds of simple comparisons captured here. Our coverage results suggest this is valid for a significant range of PHP applications.

## 6 EVALUATION

We experimentally measured the effectiveness of Apollo by using it to find faults in PHP Web applications. We designed experiments to answer the following research questions:

Q1. How many faults can Apollo find, and of what varieties?

Q2. How effective is the fault detection technique of Apollo compared to alternative approaches in terms of the number and severity of discovered faults and the line coverage achieved?

Q3. How effective is our minimization technique in reducing the size of input parameter constraints and failure-inducing inputs?

For the evaluation, we selected six open-source PHP programs from http://sourceforge.net (see Fig. 9):

- **faqforge**: tool for creating and managing documents;
- **webchess**: online chess game;

15. http://choco-solver.net/index.php?title=Main_Page.
16. Floating-point values can be handled in the same way, though none of the examined programs required it.

| program | version | #files | PHP LOC | #downloads |
|---|---|---|---|---|
| **faqforge** | 1.3.2 | 19 | 734 | 14,164 |
| **webchess** | 0.9.0 | 24 | 2,226 | 32,352 |
| **schoolmate** | 1.5.4 | 63 | 4,263 | 4,466 |
| **phpsysinfo** | 2.5.3 | 73 | 7,745 | 492,217 |
| **timeclock** | 1.0.3 | 62 | 13,879 | 23,708 |
| **phpBB2** | 2.0.21 | 78 | 16,993 | 18,668,112 |

Fig. 9. Characteristics of subject programs. The **#files** column lists the number of `.php` and `.inc` files in the program. The **PHP LOC** column lists the number of lines that contain executable PHP code. The **#downloads** column lists the number of downloads from http://sourceforge.net.

- **schoolmate**: PHP/MySQL solution for administering elementary, middle, and high schools;
- **phpsysinfo**: displays system information, e.g., uptime, CPU, memory, etc.;
- **timeclock** is a Web-based timeclock system;
- **phpBB2** is a discussion forum.

## 6.1 Generation Strategies

We use the following test input-generation strategies in the remainder of this section:

- **Apollo** generates test inputs using the technique described in Section 3.
- **Randomized** is a test input generation strategy that generates test inputs by giving random values to parameters. These values are chosen from constants harvested from the program's source code and from default values.[17] A difficulty is that the names and types of parameters are not declared explicitly in PHP applications. The Randomized algorithm discovers the names of parameters by monitoring what string values are used to access the global arrays $\_POST, $\_GET, and $\_REQUEST that are used to read the parameters supplied to PHP programs. The algorithm starts by executing the application on the empty input, and in each subsequent execution, randomly chosen harvested constants and default values are assigned to a randomly selected subset of parameters that have been detected in previous executions.[18]

To avoid bias, we ran both strategies inside the same experimental harness. This includes the Database Manager (Section 5), which supplies usernames and passwords for database access, and the UI option analyzer.

## 6.2 Methodology

To answer the first research question (Q1) we applied Apollo to the six subject programs and we classified the discovered failures into five groups based on their different failure characteristics:

- **Execution crash:** The PHP interpreter terminates with an exception.

17. Halfond and Orso [20] presented a similar static analysis for Web applications written in Java; however, they discover an approximation of the set of values for each parameter.

18. Two additional features of the Randomized algorithm are: 1) The user may supply specific values to use for specific parameters (such as username/password), and 2) it detects additional scripts that may be invoked from the output of each execution (user-input simulation).

- **Execution error:** The PHP interpreter emits an error message that is visible in the generated HTML.
- **Execution warning:** The PHP interpreter emits an error message that is invisible in the generated HTML.
- **HTML error:** The program generates HTML for which the validator produces an error report.
- **HTML warning:** The program generates HTML for which the validator produces a warning report.

This classification is a refinement of the one presented in Section 2.3.

To answer the second research question (Q2) we compared both the coverage achieved and the number of faults found with the **Randomized** generation strategy. Coverage was measured using the line coverage metric, i.e., the ratio of the number of executed lines to the total number of lines with executable PHP code in each application.

We ran each test input-generation strategy for 20 minutes on each subject program, by which time the coverage achieved and number of bugs found by Apollo starts to level off. During this time each strategy generated hundreds of inputs.

This time budget includes all experimental tasks, i.e., program execution, harvesting of constant values from program source, test generation, constraint solving (where applicable), output validation via an oracle, and line coverage measurement. For our experiments, we use the WDG offline HTML validator, version 1.2.2.

We also compared Apollo's results to the results reported by Minamide's static analysis [34] on the four subject programs that we have in common (Section 6.3.1 presents the results).

To answer the third research question, about the effectiveness of the input minimization, we performed the following experiments. Recall that several execution paths and inputs may expose the same failure. Our input minimization algorithm attempts to produce the shortest possible input that exposes each failure. The inputs to the minimizer are the failure found by the algorithm of Fig. 6 along with all the execution paths that expose each failure.

## 6.3 Results

The graphs in Figs. 11 and 12 visualize how coverage and the number of failures found increases over time, when both techniques are given up to 20 minutes. For each of the failures reported by the two algorithms, we manually investigated the problem and identified the faulty statements that caused the problem by fixing the problems and making sure that the failures did not recur. Fig. 10 tabulates the final line coverage and faults when both techniques are given 20 minutes. To the best of our knowledge, all reported faults are counted only once.

In most cases, Apollo finds many more failures and achieves much better coverage than Randomized. The only notable exception to this is **phpsysinfo** (Fig. 12b), which is not a Web application, but a program that displays system information, and therefore does not rely much on user input. In this case, using the UI analyzer, Randomized was able to perform almost as well as Apollo.

On most applications, Apollo's coverage converges after about 10 minutes while Randomized converges sooner, and at a much lower level. Apollo coverage of **webchess** (Fig. 11d)

| program | strategy | #inputs generated | line coverage % | execution | | | HTML validation | | Total faults |
|---|---|---|---|---|---|---|---|---|---|
| | | | | crash | error | warning | error | warning | |
| faqforge | Randomized | 1941 | 31.6 | 0 | 1 | 0 | 32 | 4 | 37 |
| | Apollo | 1573 | 96.5 | 0 | 6 | 0 | 66 | 21 | 93 |
| webchess | Randomized | 1805 | 5.9 | 1 | 13 | 0 | 3 | 0 | 17 |
| | Apollo | 1508 | 40.6 | 1 | 19 | 0 | 11 | 0 | 31 |
| schoolmate | Randomized | 1396 | 8.7 | 1 | 0 | 0 | 20 | 0 | 21 |
| | Apollo | 1355 | 65.6 | 3 | 25 | 0 | 95 | 0 | 123 |
| phpsysinfo | Randomized | 490 | 55.6 | 0 | 4 | 2 | 2 | 0 | 8 |
| | Apollo | 356 | 56.2 | 0 | 5 | 2 | 2 | 0 | 9 |
| timeclock | Randomized | 1089 | 5.8 | 0 | 1 | 1 | 76 | 1 | 79 |
| | Apollo | 1112 | 26.8 | 0 | 2 | 1 | 375 | 9 | 387 |
| phpbb2 | Randomized | 2497 | 11.4 | 0 | 3 | 0 | 17 | 0 | 20 |
| | Apollo | 981 | 31.7 | 0 | 3 | 5 | 22 | 0 | 30 |
| Total | Randomized | 9218 | 19.8 | 2 | 22 | 3 | 150 | 5 | 182 |
| | Apollo | 6885 | 52.9 | 4 | 60 | 8 | 571 | 30 | 673 |

Fig. 10. Experimental results for 20 minute test generation runs. The table presents results for each subject program, and each strategy, separately. The **#inputs** column presents the number of inputs that each strategy created in the given time budget. The **coverage** column lists the line coverage achieved by the generated inputs. The **execution crashes**, **errors**, and **warnings** and **HTML errors** and **warnings** columns list the number of faults in the respective categories. The **Total faults** columns sums up the number of discovered faults.
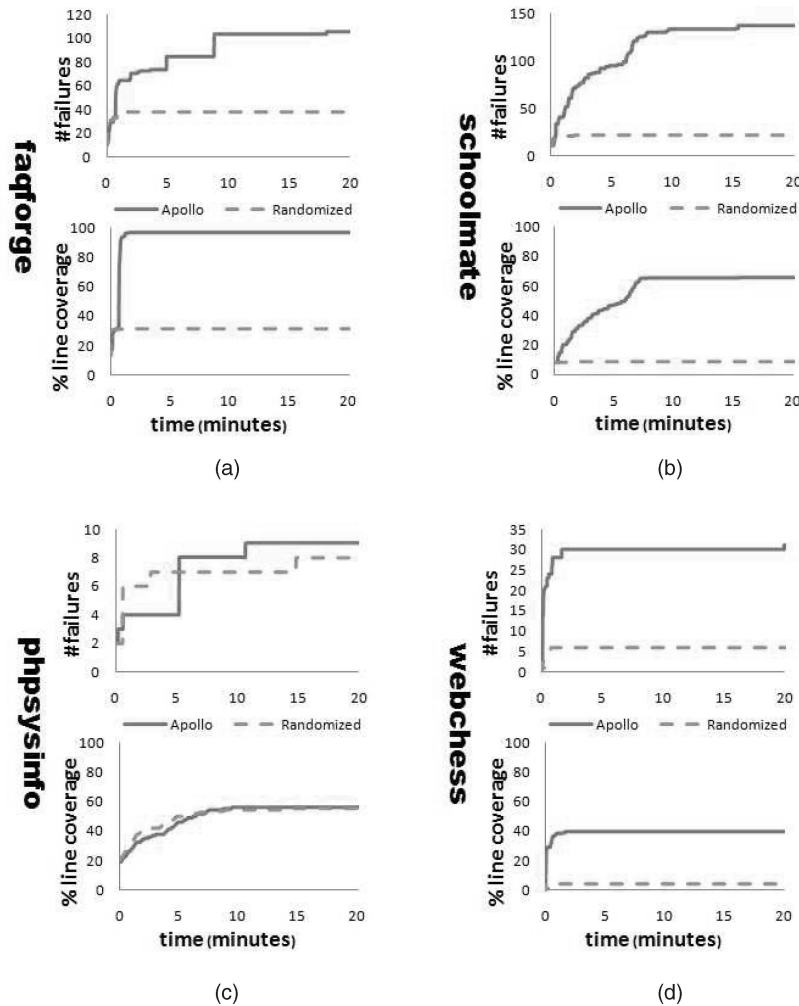


Fig. 11. Percentage of line coverage (bottom) and failures found (top) achieved by Apollo (full line) and Randomized (dashed line) input-generation techniques in 20 minutes of execution for subject programs (a) faqforge, (b) schoolmate, (c) phpsysinfo, and (d) webchess.

converges earlier since Apollo fails to login two players, and thus is unable to simulate a chess game. On **faqforge** (Fig. 11a), Apollo quickly covers almost 100 percent of the code. In both cases, new failures are still found in the covered code until the last minute of execution. We note that our results confirm the common belief that there is a very strong correlation between coverage and number of failures found.

The **Apollo** strategy found 673 faults in the subject applications versus only 182 faults for **Randomized**. Moreover, the **Apollo** test generation strategy achieved an average line coverage of 52.9 percent versus only 19.8 percent for **Randomized**.

The coverage of **phpbb2** and **timeclock** is relatively low as the output of these applications contains client-side
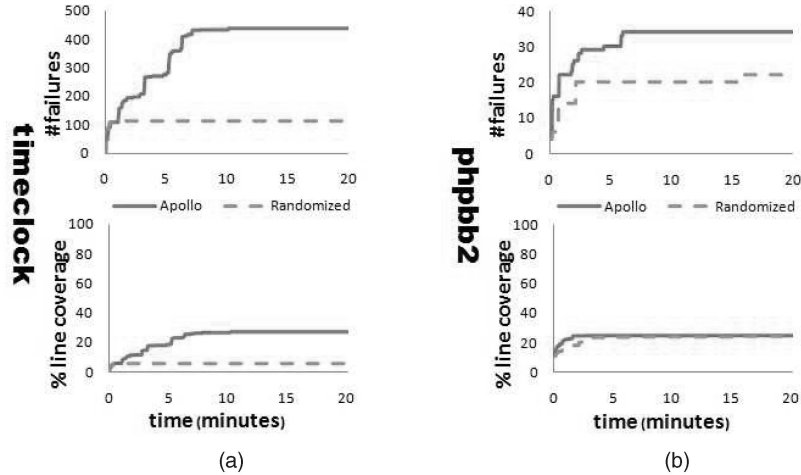
Fig. 12. Percentage of line coverage (bottom) and failures found (top) achieved by Apollo (full line) and Randomized (dashed line) input-generation techniques in 20 minutes of execution for subject programs (a) timeclock and (b) phpbb2.

scripts written in JavaScript, which Apollo currently does not analyze.

Figs. 13 and 14 classify the faults reported by Apollo. The execution errors (Fig. 13) are dominated by database-related errors, where the application had difficulties accessing the database, resulting in error messages such as 1) "supplied argument is not a valid MySQL result resource" and 2) "Unable to jump to row 0 on MySQL result." The two SQL-related error messages quoted above occurred in faqforge (9 cases of error 1) and webchess (19 cases of error 1 and 1 case of error 2), schoolmate (20 cases of error 1 and 9 cases of error 2), timeclock (1 case of error 1), and phpbb2 (1 case of error 1).

These failures have the same cause: User-supplied input parameters are concatenated directly into SQL query strings and leaving these parameters blank results in malformed SQL, which in turn causes the `mysql_query` functions to return an invalid result. The subject programs failed to check the return value of `mysql_query` and simply assume that a valid result is returned. These faults are indications of a potentially serious problem: The concatenation of user-supplied strings into SQL queries makes these programs vulnerable to SQL-injection attacks [12]. Thus, our testing approach indicates possible SQL-injection vulnerabilities despite not being specifically designed to look for security issues.

The three execution crashes (when the interpreter terminates with an exception) in Fig. 10 happen when the interpreter tries to load files or functions that are missing.

For instance, for some inputs that can be supplied to the schoolmate subject program, the PHP interpreter attempts to load a file that does not exist in the current distribution of schoolmate. Since schoolmate has 63 files and PHP is an interpreted language that allows the use of runtime string values when loading files, it is hard to detect such faults. Apollo also discovers a severe fault in the webchess subject program. This fault occurs when the interpreter tries to call to a function that is undefined since the PHP file implementing it is not included due to a value supplied as one of the parameters.

The 673 malformed HTML faults can be divided into several categories (Fig. 14). These faults are mainly concerned with HTML elements that occur in the wrong place, HTML elements with incorrect values, and with unclosed tags. The breakdown of HTML faults is similar across the different PHP applications.

| Fault Category | Faults | Percentage |
|---|---|---|
| Missing attribute | 144 | 23.7 |
| Element not allowed | 132 | 21.7 |
| No attribute | 78 | 12.8 |
| Missing end tag | 57 | 9.4 |
| Can't generate system identifier | 52 | 8.5 |
| Incorrect end tag | 35 | 5.8 |
| Undefined element | 24 | 3.9 |
| character not allowed | 23 | 3.8 |
| End tag for unfinished element | 20 | 3.3 |
| Unfinished tag | 19 | 3.1 |
| Incorrect attribute | 6 | 1.0 |
| Duplicate specification | 5 | 0.8 |
| Invalid comment declaration | 3 | 0.5 |
| Incorrect attribute value | 3 | 0.5 |
| Not a member of a group | 2 | 0.3 |
| Start tag omitted | 1 | 0.2 |
| Prolog can't be omitted | 1 | 0.2 |
| Declaration not allowed | 1 | 0.2 |
| Attribute must start with name | 1 | 0.2 |

| Fault Category | Faults | Percentage |
|---|---|---|
| Malformed SQL | 48 | 66.7 |
| Incorrect parameter | 8 | 11.1 |
| Resource used as offset | 4 | 5.6 |
| Failed to open stream | 4 | 5.6 |
| File not found | 2 | 2.7 |
| Can't open connection | 2 | 2.7 |
| Assigning reference | 2 | 2.7 |
| Undefined function | 2 | 2.7 |

Fig. 13. Classification of the execution faults found by Apollo.

Fig. 14. Classification of the HTML faults found by Apollo.

| program | success rate% | path constraints | | inputs | |
|---|---|---|---|---|---|
| | | orig. size | reduction | orig. size | reduction |
| faqforge | 64 | 22.3 | 78% | 9.3 | 69% |
| webchess | 91 | 23.4 | 81% | 10.9 | 60% |
| schoolmate | 51 | 22.9 | 62% | 11.5 | 42% |
| phpsysinfo | 82 | 24.3 | 82% | 17.5 | 74% |
| timeclock | 54 | 19.3 | 52% | 8.9 | 39% |
| phpBB2 | 72 | 25.5 | 59% | 14.2 | 41% |

Fig. 15. Results of minimization. The **success rate** indicates the percentage of faults whose exposing input was successfully minimized (i.e., the minimizer found a shorter exposing input). The **orig. size** columns list the average size of original (un-minimized) path constraints and inputs. The size of a path constraint is the number of conjuncts. The size of an input is the number of key-value pairs in the input. The **reduction** columns list the amount by which the minimized size is smaller than the unminimized size (i.e., $1 - \frac{minimized}{unminimized}$). The higher the percentage, the more often minimization is helpful.

### 6.3.1 Comparison with Static Analysis

Minamide [34] presents a static analysis for discovering HTML malformedness faults in PHP applications. Minamide's analysis tool approximates the string output of a program with a context-free grammar, and then discovers unmatched tags by intersecting this grammar with the regular expression of matched pairs of delimiters (open/closed tags). In contrast, our analysis uses an HTML validator and covers the entire language standard.

We performed our evaluation on four of the benchmarks studied by Minamide (webchess, faqforge, schoolmate, and timeclock), and we compare the number of faults found by both tools (using a time budget of 20 minutes for Apollo). For these four subject programs, Apollo found 2.2 times as many unmatched tags failures as Minamide (125[19] versus 56). The faults found by Minamide's tool are not publicly available so we do not know whether Apollo discovered all faults that Minamide's tool discovered. However, Apollo also found 453 other malformed HTML failures and 66 execution faults, both of which are out of reach for Minamide's tool.

### 6.3.2 Path Constraint Minimization

We measure the effectiveness of the minimization algorithm of Section 3.3 via the reduction ratio between the size of the shortest original (unminimized) path constraint (and input) and the minimized path constraint (and input).

Fig. 15 tabulates the results. The results show that our input minimization technique effectively reduces the size of inputs by at least 42 percent. The minimization technique managed to reduce the exposing input size for more than 50 percent of the faults.

## 6.4 Threats to Validity

**Construct validity.** Why do we count malformed HTML as a defect in dynamically generated webpages? Does a webpage with malformed HTML pose a real problem or this is an artificial problem generated by the overly conservative specification of the HTML language? Although Web browsers are resilient to malformed HTML, we have encountered cases when malformed HTML crashed the popular Internet Explorer Web browser. More importantly, even though a particular browser might

19. 125 is the sum of Malformed HTML failures related to tags: Missing end tag, unopened close tag, end tag for unfinished element, and unfinished tag.

tolerate malformed HTML, different browsers or different versions of the same browser may not display all information in the presence of malformed HTML. This becomes crucial for some websites; for example, for sites related to financial transactions. Many websites provide a button for verifying the validity of statically generated HTML. The challenges of dynamically generated webpages prevent the same institutions from validating the content.

Why do we use line coverage as a quality metric? We use line coverage only as a *secondary* metric, our *primary* metric being the number of faults found. Line coverage indicates how much of the application was explored by the analysis. An analysis can only find faults in lines that are covered, so more coverage generally leads to more faults being detected.

Why do we present the user with minimized path constraints and inputs in addition to the inputs exposing the failure? Although an input that corresponds to a longer path constraint still exposes the same failure, in our experience, the removal of superfluous information helps programmers with pinpointing the location of the fault.

**Internal validity.** Did Apollo discover real, unseeded, and unknown faults? Since we used subject projects developed by others, we could not influence the quality of the subject programs. Apollo does not search for known or seeded faults, but it finds *real* faults in real programs. For those subject programs that connect to a database, we populated the database with random records. The only thing that is "seeded" into the experiment is a username/ password combination, so that Apollo can access the records stored in the database.

**External validity.** Will our results generalize beyond the subject programs? We only used Apollo to find faults in six PHP projects. These may have serious quality problems, or may be unrepresentative in other ways. Four of the subject programs were also used as subject programs by Minamide [34]. We chose two additional subject programs, phpsysinfo and phpBB2, in order to ensure we use high quality programs. Both programs were chosen using the following three criteria: applications with more than 10K LOC, that had been edited recently (i.e., recent CVS commits), and have a large number of downloads (more than 100,000). Both programs are ranked among the top 0.5 percent projects on sourceforge (phpBB is ranked the 29th most popular project out of all 154,880 sourceforge applications, as of 27 August 2009). Nevertheless, Apollo managed to find faults in both of these programs.

While it is true that Apollo finds hundreds of faults in the subject programs, this does not mean that these applications are of poor quality. Many of the faults Apollo discovers are malformed HTML faults. PHP applications tend to have many malformed HTML faults for several reasons. First, it is hard to avoid malformed HTML faults when writing a program in one language that generates code in another. Second, for the same reason, it is hard to find malformed HTML faults with manual inspection and code reviews. Third, in order to find malformed HTML faults automatically, the developer needs to generate and inspect all possible output of the applications. Fourth, malformed HTML faults usually have low priority when fixing faults as browsers are attempting to recover from them automatically.

Similarly, the fact the randomized version finds many faults does not mean that the applications are of low

quality. Even the randomized version is doing a considerable amount of work to automate the discovery of faults. It discovers input parameters, constants, and simulates user input, in order to find as many faults as it does.

**Reliability.** Are the results reproducible? The subject programs that we used are publicly available from sourceforge. The faults that we found are available for examination at http://pag.csail.mit.edu/apollo.

## 6.5 Limitations

**Simulating user inputs based on locally executed JavaScript**. The HTML output of a PHP script might contain buttons and arbitrary snippets of JavaScript code that are executed when the user presses the corresponding button. The actions that the JavaScript interpreter might perform are currently not analyzed by Apollo. For instance, the JavaScript code might pass specific arguments to the PHP script. As a result, Apollo might report false positives. For example, Apollo might report a false positive if Apollo decides to execute a PHP script as a result of simulating a user pressing a button that is not visible. Apollo might also report a false positive if it attempts to set an input parameter that would have been set by the JavaScript code. In our experiments, Apollo did not report any false positives.

**Limited tracking in native methods.** Apollo has limited tracking of input parameters through PHP native methods. PHP native methods are implemented in C, which make it difficult to automatically track how input parameters are transformed into output parameters. We have modified the PHP interpreter to track parameters across a very small subset of the PHP native methods. Similarly to [44], we plan to create an external language to model the dependencies between inputs and outputs for native methods to increase Apollo line coverage when native methods are executed.

**Limited tracking of input parameters through the database**. Apollo does not track input parameters through the database. Thus, Apollo might not be able to explore call sequences in which subsequent calls depend on specific values of input parameters stored in the database by earlier calls. It is possible to extend Apollo to track input parameters through the database in the same way Emmi et al. [14] extended concolic testing to database applications.

**Limited sources of input parameters.** Apollo currently considers parameters as only inputs coming from the global arrays _POST, _GET, and _REQUEST. Supporting other global parameters such as _ENV and _COOKIE is straightforward.

**Limited forms of constraints to be solved.** In theory, Apollo might be unable to cover certain parts of an application because the constraints that use are fairly simple. However, from what we have observed in our experiments so far, PHP Web applications do not seem to manipulate their input parameters in complex ways, and the very simple constraints that can be solved with Choco have been adequate for our purposes. We have only observed a very few isolated cases where the solving of more complex constraints (see, e.g., [28]) would have helped.

## 7  RELATED WORK

An earlier version of this paper was presented at ISSTA '08 [2]. The Apollo tool presented there did not handle the problem of automatically simulating user interactions in Web applications. Instead, it relied on a manual transformation of the program under test to enable the exploration of a few selected user inputs. The current paper also extends [2] by providing a more extensive evaluation, which includes two new large Web applications, and by presenting a detailed classification of the faults found by Apollo. In addition, the Apollo tool presented in [2] did not yet support Web server integration.

In the remainder of this section, we discuss three categories of related work: 1) combined concrete and symbolic execution, 2) techniques for input minimization, and 3) testing of Web applications.

## 7.1  Combined Concrete and Symbolic Execution

DART [17] is a tool for finding combinations of input values and environment settings for C programs that trigger errors such as assertion failures, crashes, and nontermination. DART combines random test generation with symbolic reasoning to keep track of constraints for executed control-flow paths. A constraint solver directs subsequent executions toward uncovered branches. Experimental results indicate that DART is highly effective at finding large numbers of faults in several C applications and frameworks, including important and previously unknown security vulnerabilities. CUTE [39] is a variation (called *concolic testing*) on the DART approach. The authors of CUTE introduce a notion of approximate pointer constraints to enable reasoning over memory graphs and handle programs that use pointer arithmetic.

Subsequent work extends the original approach of combining concrete and symbolic executions to accomplish two primary goals: 1) improving scalability [1], [5], [15], [16], [18], [32], and 2) improving execution coverage and fault detection capability through better support for pointers and arrays [7], [39], better search heuristics [18], [25], [31], or by encompassing wider domains such as database applications [14].

Godefroid [15] proposed a compositional approach to improve the scalability of DART. In this approach, summaries of lower level functions are computed dynamically when these functions are first encountered. The summaries are expressed as pre and postconditions of the function in terms of its inputs. Subsequent invocations of these lower level functions reuse the summary. Anand et al. [1] extend this compositional approach to be demand-driven to reduce the summary computation effort.

Exploiting the structure of the program input may improve scalability [16], [32]. Majumdar and Xu [32] abstract context-free grammars that represent the program inputs to produce a symbolic grammar. This grammar reduces the number of input strings to enumerate during test generation.

Majumdar and Sen [31] describe hybrid concolic testing, which interleaves random testing with bounded exhaustive symbolic exploration to achieve better coverage. Inkumsah and Xie [25] combine evolutionary testing using genetic mutations with concolic testing to produce longer sequences of test inputs. SAGE [18] also uses improved heuristics, called *white box fuzzing*, to achieve higher branch coverage.

Emmi et al. [14] extend concolic testing to database applications. This approach creates and inserts database

records and enables testing program code that depends on embedded SQL queries.

Wassermann et al. present a concolic-testing tool [45] for identifying SQL-injection vulnerabilities in PHP applications by using dynamic input generation and a string analysis [34]. Our work is related to Wassermann et al.'s but differs significantly in goals and technique. We discuss the similarities (1-3) and differences (4-7) below.

1. Both tools aim at creating inputs that expose faults in PHP programs and at achieving high code coverage.
2. Both tools use dynamic analysis to track the flow of input data though the execution of the program, collect symbolic constraints, and use a constraint solver to create additional inputs.
3. Both tools use test oracles to check the results of the execution.
4. The focus of the work of Wassermann et al. is software security, while the focus of ours is correctness of the output. Their tool aims to construct inputs that expose known SQL-injection vulnerabilities (the tool requires an indication of a vulnerability). In contrast, our tool aims to construct inputs that lead to previously unknown faults that result in invalid output.
5. Their tool is not fully automatic, while ours is. Their tool requires manual loading of pages and supplying of inputs to the page. Also, it is unable to explore realistic executions which result from a sequence of user interactions with the program. In contrast, our tool automatically handles interactive user input, including discovering new inputs by analysis of generated HTML pages and propagation of symbolic inputs through sessions and cookies. Full automatization of user interaction is a major contribution of this work.
6. Their tool requires a specialized solver for a logic of string constraints that includes strings, integers, regular-language constraints, and array constraints. Their constraint solver is based on finite-state automata and transducers. The authors do not discuss the algorithmic complexity of the solving problem or the empirical efficiency of the solver. In contrast, our tool creates constraints that can be efficiently solved by any linear-arithmetic solver. In principle, a potential advantage of using a more complex constraint theory is that it may allow achieving higher coverage. In practice, however, using a simple but efficient constraint theory works well. Our tool achieved high line coverage, 52.9 percent (Wassermann et al. mention that the aim of their work is to "achieve a designated code-coverage metric" but provide no coverage results for their experiments).
7. Their approach to concolic execution relies on performing source-code instrumentation and backward-slice computation by reexecuting the same input multiple times and instrumenting additional code.

As the authors mention, this "is effective only when the points of possible failure are known and relatively localized." In contrast, our tool works on unchanged application code, generates symbolic constraints directly using a modified interpreter, and does execute any input more than once. These features make our approach to concolic execution more generally applicable.

Some approaches aim at checking functional correctness. A number of tools [4], [6] use a separate implementation of the function being tested to compare outputs. This limits the approach to situations where a second implementation exists.

While our work builds on this significant body of research, there are two significant differences. First, our work goes beyond simple assertion failures and crashes by using on an oracle (in the form of an HTML validator) to determine correctness, which means that our tool can handle situations where the program has functionally incorrect behavior without relying on programmer assertions. Second, our work addresses PHP's complex execution model, which involves multiple scripts invoked via user-interface options in generated HTML pages, and communicating values via session state and cookies. The only other concolic-testing approach for PHP [45] does not present a fully automatic solution for dealing with multiple interrelated PHP scripts.

## 7.2 Minimizing Failure-Inducing Inputs

Our work minimizes the constraints on the input parameters. This shortens the failure-inducing inputs and helps with pinpointing the cause of faults. Godefroid et al. [18] faced this challenge since their technique produces several distinct inputs that expose the same fault. Their approach hashes all such inputs and returns an example failure-inducing input. Our work also addresses another issue: identifying the minimal set of program variables that are essential to induce the failure. In this regard, our work is similar to *delta debugging* [9], [47] and its extension *hierarchical delta debugging* [35]. These approaches modify the failure-inducing input directly, thus leading to a single, minimal failure-inducing input. In contrast, our technique modifies the set of constraints on the failure-inducing input. This creates minimal *patterns* of failure-inducing inputs, which facilitates debugging. Moreover, our technique is more efficient because it takes advantage of the (partial) overlapping of different inputs.

Two recent papers present other approaches to minimizing faulty input. Clause and Orso [8] present a technique, based on dynamic tainting, to find the minimal subset of input responsible for the fault. Given a program that fails with a runtime exception caused by an incorrect value, Sinha et al.'s technique [40] identifies the source statement at which the incorrect assignment was made, which helps to locate the fault.

## 7.3 Testing of Web Applications

Existing techniques for fault detection in Web applications focus on output correctness and security.

Minamide [34] uses static string analysis and language transducers to model PHP string operations to generate *potential* HTML output—represented by a context-free grammar—from the Web application. This method can be used to generate HTML document instances of the resulting grammar and to validate them using an existing HTML validator. As a more complete alternative, Minamide

proposes a *matching validation* which checks for containment of the generated context-free grammar against a regular subset of the HTML specification. However, this approach can only check for matching start and end tags in the HTML output, while our technique covers the entire HTML specification. Also, flow-insensitive and context-insensitive approximations in the static analysis techniques used in this method result in false positives, while our method reports only real faults.

Benedikt et al. [3] present a tool, VeriWeb, for automatically testing dynamic webpages. They use a model checker to systematically explore all paths (up to a certain bound) of user navigatable components in a website. When the exploration encounters HTML forms, VeriWeb uses *SmartProfiles*. SmartProfiles are user-specified attribute-value pairs that are used to automatically populate forms and supply values that should be provided as inputs. Although VeriWeb can automatically fill in the forms, the human tester needs to prepopulate the user profiles with values that a user would provide. Similarly, the WAVES tool by Huang et al. [23] performs automatic form completion by using a textual analysis to associate "topics" with input values that occur in HTML forms, in combination with a self-learning knowledge base that associates values with topics. In contrast, Apollo automatically discovers input values based on the examination of branch conditions on execution paths. Benedikt et al. do not report any faults found, while we report 673.

Halfond and Orso [20] use static analysis of the server-side implementation logic to extract a Web application's interface, i.e., the set of input parameters and their potential values. (Halfond et al. later extended that approach [19] to include dynamic analysis and specialized constraint solving.) They implemented their technique for Web applications written in Java. They obtained better code coverage with test cases based on the interface extracted using their technique as compared to the test cases based on the interface extracted using a conventional Web crawler. However, the coverage may depend on the choices made by the test generator to combine parameter values—an exhaustive combination of values may be needed to maximize code coverage. In contrast, our work uses dynamic analysis of server-side implementation logic for fault detection and minimizes the number of inputs needed to maximize the coverage. Furthermore, we include results on fault detection capabilities of our technique.

Like us, McAllister et al. [33] tackle the problem of testing interactive Web applications. Their method relies on prerecorded traces of user interactions, while our approach automatically discovers allowable interactions. Moreover, their approach to handling persistent state relies on instrumenting one particular Web application framework, Django. In contrast, our approach is to instrument the PHP runtime system and observe database interactions. This allows handling state of PHP applications regardless of any framework they may use.

Dynamic analysis of string values generated by Web applications has been considered in a *reactive* mode to prevent the execution of insidious commands (*intrusion prevention*) and to raise an alert (*intrusion detection*) [26], [30],

[36], [38], [41]. As far as we know, our work is the first attempt at *proactive* fault detection in Web applications using dynamic analysis.

Existing techniques for detection of security vulnerabilities (e.g., *SQL injection* or *cross-site scripting*) in Web applications use static analysis [24], [27], [34], [43], [46] or dynamic analysis [29].

Kieżun et al. present a dynamic tool, Ardilla [29], to create SQL and XSS attacks. Their tool uses dynamic tainting, concolic execution, and attack-candidate generation and validation. Like ours, their tool reports only real faults. However, Kieżun et al. focus on finding security faults, while we concentrate on functional correctness. Their tool builds on and extends the input-generation component of Apollo, but does not address the problem of user interaction. It is an area of future research to combine Apollo's user interaction and state matching with Ardilla's exploit-detection capabilities.

## 8 CONCLUSIONS

We have presented a technique for finding faults in PHP Web applications that is based on combined concrete and symbolic execution. The work is novel in several respects. First, the technique not only detects runtime errors but also uses an HTML validator as an oracle to determine situations where malformed HTML is created. Second, we address a number of PHP-specific issues, such as the simulation of interactive user input that occurs when user-interface elements on generated HTML pages are activated, resulting in the execution of additional PHP scripts. Third, we perform an automated analysis to minimize the size of failure-inducing inputs.

We created a tool, Apollo, that implements the analysis. We evaluated Apollo on six open-source PHP web applications. Apollo's test generation strategy achieves over 50 percent line coverage. Apollo found a total of 673 faults in these applications: 72 execution problems and 601 cases of malformed HTML. Finally, Apollo also minimizes the size of failure-inducing inputs: The minimized inputs are up to $5.3\times$ smaller than the unminimized ones.

## REFERENCES

[1] S. Anand, P. Godefroid, and N. Tillmann, "Demand-Driven Compositional Symbolic Execution," *Proc. Int'l Conf. Tools and Algorithms for the Construction and Analysis of Systems,* pp. 367-381, 2008.

[2] S. Artzi, A. Kiezun, J. Dolby, F. Tip, D. Dig, A. Paradkar, and M.D. Ernst, "Finding Bugs in Dynamic Web Applications," *Proc. Int'l Symp. Software Testing and Analysis,* pp. 261-272, 2008.

[3] M. Benedikt, J. Freire, and P. Godefroid, "VeriWeb: Automatically Testing Dynamic Web Sites," *Proc. Int'l Conf. World Wide Web,* 2002.

[4] D. Brumley, J. Caballero, Z. Liang, J. Newsome, and D. Song, "Towards Automatic Discovery of Deviations in Binary Implementations with Applications to Error Detection and Fingerprint Generation," *Proc. 16th USENIX Security Symp.,* 2007.

[5] C. Cadar, D. Dunbar, and D.R. Engler, "Klee: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs," *Proc. USENIX Symp. Operating Systems Design and Implementation,* pp. 209-224, 2008.

[6] C. Cadar and D.R. Engler, "Execution Generated Test Cases: How to Make Systems Code Crash Itself," *Proc. Int'l SPIN Workshop Model Checking of Software,* pp. 2-23, 2005.

[7] C. Cadar, V. Ganesh, P.M. Pawlowski, D.L. Dill, and D.R. Engler, "EXE: Automatically Generating Inputs of Death," *Proc. Conf. Computer and Comm. Security*, pp. 322-335, 2006.

[8] J. Clause and A. Orso, "Penumbra: Automatically Identifying Failure-Relevant Inputs Using Dynamic Tainting," *Proc. Int'l Symp. Software Testing and Analysis*, 2009.

[9] H. Cleve and A. Zeller, "Locating Causes of Program Failures," *Proc. Int'l Conf. Software Eng.*, pp. 342-351, 2005.

[10] H. Cleve and A. Zeller, "Locating Causes of Program Failures" *Proc. Int'l Conf. Software Eng.*, pp. 342-351, May 2005.

[11] C. Csallner, N. Tillmann, and Y. Smaragdakis, "DySy: Dynamic Symbolic Execution for Invariant Inference," *Proc. Int'l Conf. Software Eng.*, pp. 281-290, 2008.

[12] D. Dean and D. Wagner, "Intrusion Detection via Static Analysis," *Proc. Symp. Research in Security and Privacy*, pp. 156-169, May 2001.

[13] C. Demartini, R. Iosif, and R. Sisto, "A Deadlock Detection Tool for Concurrent Java Programs," *Software—Practice and Experience*, vol. 29, no. 7, pp. 577-603, June 1999.

[14] M. Emmi, R. Majumdar, and K. Sen, "Dynamic Test Input Generation for Database Applications," *Proc. Int'l Symp. Software Testing and Analysis*, pp. 151-162, 2007.

[15] P. Godefroid, "Compositional Dynamic Test Generation," *Proc. Ann. Symp. Principles of Programming Languages*, pp. 47-54, 2007.

[16] P. Godefroid, A. Kieżun, and M.Y. Levin, "Grammar-Based Whitebox Fuzzing," *Proc. ACM SIGPLAN Conf. Programming Language Design and Implementation*, pp. 206-215, 2008.

[17] P. Godefroid, N. Klarlund, and K. Sen, "DART: Directed Automated Random Testing," *Proc. ACM SIGPLAN Conf. Programming Language Design and Implementation*, pp. 213-223, 2005.

[18] P. Godefroid, M.Y. Levin, and D. Molnar, "Automated Whitebox Fuzz Testing," *Proc. Network Distributed Security Symp.*, pp. 151-166, 2008.

[19] W.G. Halfond, S. Anand, and A. Orso, "Precise Interface Identification to Improve Testing and Analysis of Web Applications," *Proc. Int'l Symp. Software Testing and Analysis*, 2009.

[20] W.G.J. Halfond and A. Orso, "Improving Test Case Generation for Web Applications Using Automated Interface Discovery," *Proc. Joint Meeting European Software Eng. Conf. and ACM SIGSOFT Symp. Foundations of Software Eng.*, pp. 145-154, 2007.

[21] K. Havelund and T. Pressburger, "Model Checking Java Programs Using Java PathFinder," *Int'l J. Software Tools for Technology Transfer*, vol. 2, no. 4, pp. 366-381, 2000.

[22] G.J. Holzmann, "The Model Checker SPIN," *Software Eng.*, vol. 23, no. 5, pp. 279-295, 1997.

[23] Y.-W. Huang, S.-K. Huang, T.-P. Lin, and C.-H. Tsai, "Web Application Security Assessment by Fault Injection and Behavior Monitoring," *Proc. 12th Int'l Conf. World Wide Web*, pp. 148-159, 2003.

[24] Y.-W. Huang, F. Yu, C. Hang, C.-H. Tsai, D.T. Lee, and S.-Y. Ku, "Verifying Web Applications Using Bounded Model Checking," *Proc. Int'l Conf. Dependable Systems and Networks*, pp. 199-208, 2004.

[25] K. Inkumsah and T. Xie, "Evacon: A Framework for Integrating Evolutionary and Concolic Testing for Object-Oriented Programs," *Proc. IEEE/ACM Int'l Conf. Automated Software Eng.*, 2007.

[26] M. Johns and C. Beyerlein, "SMask: Preventing Injection Attacks in Web Applications by Approximating Automatic Data/Code Separation," *Proc. ACM Symp. Applied Computing*, 2007.

[27] N. Jovanovic, C. Kruegel, and E. Kirda, "Pixy: A Static Analysis Tool for Detecting Web Application Vulnerabilities (Short Paper)," *Proc. IEEE Symp. Security and Privacy*, pp. 258-263, 2006.

[28] A. Kieżun, V. Ganesh, P.J. Guo, P. Hooimeijer, and M.D. Ernst, "HAMPI: A Solver for String Constraints," *Proc. Int'l Symp. Software Testing and Analysis*, 2009.

[29] A. Kieżun, P. Guo, K. Jayaraman, and M. Ernst, "Automatic Creation of SQL Injection and Cross-Site Scripting Attacks," *Proc. Int'l Conf. Software Eng.*, pp. 199-209, 2009.

[30] B. Livshits, M. Martin, and M.S. Lam, "SecuriFly: Runtime Protection and Recovery from Web Application Vulnerabilities," technical report, Stanford Univ., 2006.

[31] R. Majumdar and K. Sen, "Hybrid Concolic Testing," *Proc. Int'l Conf. Software Eng.*, pp. 416-426, 2007.

[32] R. Majumdar and R.-G. Xu, "Directed Test Generation Using Symbolic Grammars," *Proc. IEEE/ACM Int'l Conf. Automated Software Eng.*, pp. 134-143, 2007.

[33] S. McAllister, E. Kirda, and C. Kruegel, "Leveraging User Interactions for In-Depth Testing of Web Applications," *Proc. 11th Int'l Symp. Recent Advances in Intrusion Detection*, pp. 191-210, 2008.

[34] Y. Minamide, "Static Approximation of Dynamically Generated Web Pages," *Proc. Int'l Conf. World Wide Web* 2005.

[35] G. Misherghi and Z. Su, "HDD: Hierarchical Delta Debugging," *Proc. Int'l Conf. Software Eng.*, pp. 142-151, 2006.

[36] A. Nguyen-Tuong, S. Guarnieri, D. Greene, J. Shirley, and D. Evans, "Automatically Hardening Web Applications Using Precise Tainting," *Proc. Int'l Conf. Information Security*, 2005.

[37] R. O'Callahan, personal communication, 2008.

[38] T. Pietraszek and C.V. Berghe, "Defending against Injection Attacks through Context-Sensitive String Evaluation," *Proc. Recent Advances in Intrusion Detection*, pp. 124-145, 2005.

[39] K. Sen, D. Marinov, and G. Agha, "CUTE: A Concolic Unit Testing Engine for C," *Proc. ACM SIGSOFT Int'l Symp. Foundations of Software Eng.*, pp. 263-272, 2005.

[40] S. Sinha, H. Shah, C. Görg, S. Jiang, and M. Kim, "Fault Localization and Repair for Java Runtime Exceptions," *Proc. Int'l Symp. Software Testing and Analysis*, 2009.

[41] Z. Su and G. Wassermann, "The Essence of Command Injection Attacks in Web Applications," *Proc. Ann. Symp. Principles of Programming Languages*, pp. 372-382, 2006.

[42] W. Visser, C.S. Păsăreanu, and R. Pelánek, "Test Input Generation for Java Containers Using State Matching," *Proc. Int'l Symp. Software Testing and Analysis*, pp. 37-48, 2006.

[43] G. Wassermann and Z. Su, "Sound and Precise Analysis of Web Applications for Injection Vulnerabilities," *Proc. ACM SIGPLAN Conf. Programming Language Design and Implementation*, pp. 32-41, 2007.

[44] G. Wassermann and Z. Su, "Static Detection of Cross-Site Scripting Vulnerabilities," *Proc. Int'l Conf. Software Eng.*, pp. 171-180, 2008.

[45] G. Wassermann, D. Yu, A. Chander, D. Dhurjati, H. Inamura, and Z. Su, "Dynamic Test Input Generation for Web Applications," *Proc. ACM/SIGSOFT Int'l Symp. Software Testing and Analysis*, pp. 249-260, 2008.

[46] Y. Xie and A. Aiken, "Static Detection of Security Vulnerabilities in Scripting Languages," *Proc. Conf. USENIX Security Symp.*, pp. 179-192, 2006.

[47] A. Zeller, "Yesterday, My Program Worked. Today, It Does Not. Why?" *Proc. ACM SIGSOFT Int'l Symp. Foundations of Software Eng.*, pp. 253-267, 1999.

[48] F. Zoufaly,, "Web Standards and Search Engine Optimization (SEO)—Does Google Care About the Quality of Your Markup?" 2008.

**Shay Artzi** received the PhD degree in computer science from the Massachusetts Institute of Technology (MIT) with a thesis entitled "Dynamically Fighting Bugs: Detection, Prevention, and Elimination," and the MS and BS degrees in computer science from the Technion (Israel Institute of Technology). He is a researcher at IBM's Thomas J. Watson Research Center in Hawthorne, New York. He has published numerous conference papers and journal articles on various aspects of improving software quality.

**Adam Kieżun** received the doctorate degree from the Massachusetts Institute of Technology, doing research on refactoring and on using automated formal-method tools to improve software testing. He is a researcher at Harvard Medical School, working on computational systems biology and human genetics. He previously worked at IBM Research in Switzerland, where he was one of the founding developers of the refactoring engine in the Eclipse Java IDE.

**Julian Dolby** has been a research staff member at the IBM Thomas J. Watson Research Center since 2000. He has worked in the past on VM technology; he currenlty works on scalable reasoning technologies for ontologies, program analysis of a range of languages, concurrent programming, and program testing.

**Frank Tip** has been at IBM Research since 1995, where he is currently managing the Program Analysis and Transformation Group. His current research interests include bug finding and fault localization, refactoring, applications of program analysis in collaborative development, and the design of data-centric mechanisms for synchronization in object-oriented programming languages.

**Danny Dig** received the PhD degree from the University of Illinois at Urbana-Champaign on automated upgrading of software applications to use the newer APIs of software libraries. After receiving the PhD, he joined the Massachusetts Institute of Technology as a postdoctoral researcher, working on interactive program transformations for retrofitting concurrency into existing sequential applications, then he returned to the University of Illinois. He is a principal investigator in the Universal Parallel Computing Research Center (UPCRC) at the University of Illinois, where he leads research on refactorings for parallelism. He started the ACM Workshop on Refactoring Tools, now in its third instance, and has served on several program committees in software engineering conferences, including the Workshop on Multicore Software Engineering and OOPSLA.

**Amit Paradkar** has been at IBM Research since 1996, where he is currently managing the Software and Services Testing Group. His current research interests include requirements elicitation techniques, application of natural language analysis to the domain of software requirements, holistic requirements analysis where requirements are defined using multiple techniques—such as text, sketches, and diagrams, and requirements-based test generation. He is a senior member of the IEEE.

**Michael D. Ernst** is an associate professor in the Computer Science and Engineering Department at the University of Washington. His research aims to make software more reliable, more secure, and easier (and more fun!) to produce. His primary technical interests are in software engineering and related areas, including programming languages, type theory, security, program analysis, bug prediction, testing, and verification. His research combines strong theoretical foundations with realistic experimentation, with an eye to changing the way that software developers work. He was previously a tenured professor at the Massachusetts Institute of Technology, and before that was a researcher at Microsoft Research.

▷ **For more information on this or any other computing topic, please visit our Digital Library at** www.computer.org/publications/dlib.