

# Merlin: Specification Inference for Explicit Information Flow Problems

Anindya Banerjee  
KSU

Benjamin Livshits, Aditya V. Nori, and Sriram K. Rajamani  
Microsoft Research

December 15, 2008

## Abstract

The last several years have seen a proliferation of static and runtime analysis tools for finding security violations that are caused by *explicit information flow* in programs. Much of this interest has been caused by the increase in the number of vulnerabilities such as cross-site scripting and SQL injections. In fact, these explicit information flow vulnerabilities commonly found in Web applications now outnumber vulnerabilities such as buffer overruns common in type-unsafe languages such as C and C++. Tools checking for these vulnerabilities require a specification to operate. In most cases the task of providing such a specification is delegated to the user. Moreover, the efficacy of these tools is only as good as the specification. Unfortunately, writing a comprehensive specification presents a major challenge: parts of the specification are easy to miss leading to missed vulnerabilities; similarly, incorrect specifications may lead to false positives.

This paper proposes MERLIN, a new algorithm for automatically inferring explicit information flow specifications from program code. Such specifications greatly reduce manual labor, and enhance the quality of results, while using tools that check for security violations caused by explicit information flow. Beginning with a data propagation graph, which represents interprocedural flow of information in the program, MERLIN aims to automatically infer an information flow specification. MERLIN models information flow paths in the propagation graph using probabilistic constraints. A naïve modeling requires an exponential number of constraints, one per path in the propagation graph. For scalability, we approximate these path constraints using constraints on chosen triples of nodes, resulting in a cubic number of constraints. We characterize this approximation as a probabilistic abstraction, using the theory of probabilistic refinement developed by McIver and Morgan. We solve the resulting system of probabilistic constraints using factor graphs, which are a well-known structure for performing probabilistic inference.

We experimentally validate the MERLIN approach by applying it to 10 large business-critical Web applications that have been analyzed with CAT.NET, a state-of-the-art static analysis tool for .NET. We find a total of 167 new confirmed specifications, which result in a total of 302 *additional* vulnerabilities across the 10 benchmarks. More accurate specifications also reduce the false positive rate: in our experiments, MERLIN-inferred specifications result in 13 false positives being removed; this constitutes a 15%

reduction in the CAT.NET false positive rate on these 10 programs. The final false positive rate for CAT.NET *after* applying MERLIN in our experiments drops to under 1%.

# 1 Introduction

Constraining information flow is fundamental to security: we do not want secret information to reach untrusted principals (confidentiality), and we do not want untrusted principals to corrupt trusted information (integrity). If we take confidentiality and integrity to the extreme, then principals from different levels of trust can never interact, and the resulting system becomes unusable. For instance, such a system would never allow a trusted user to view untrusted content from the Internet.

Thus, practical systems compromise on such extremes, and allow flow of *sanitized* information across trust boundaries. For instance, it is unacceptable to take a string from untrusted user input, and use it as part of an SQL query, since it leads to SQL injection attacks. However, it is acceptable to first pass the untrusted user input through a trusted *sanitization function*, and then use the sanitized input to construct a SQL query. Similarly, confidential data needs to be cleansed to avoid information leaks. Practical checking tools that have emerged in recent years [1, 5, 23] typically aim to ensure that all explicit flows of information across trust boundaries are sanitized.

The fundamental program abstraction used in the sequel (as well as by existing tools) is what we term the *propagation graph* — a directed graph that models all interprocedural explicit information flow in a program.<sup>1</sup> The nodes of a propagation graph are methods, and edges represent explicit information flow between methods. There is an edge from node  $m_1 \rightarrow m_2$  whenever there is a flow of information from method  $m_1$  to method  $m_2$ , through a method parameter, or through a return value, or by way of an indirect update through a pointer.

Following the widely accepted Perl taint terminology conventions [30] more precisely defined in [15], nodes of the propagation graph are classified as *sources*, *sinks*, and *sanitizers*; nodes not falling in the above categories are termed *regular* nodes. A source node returns tainted data whereas it is an error to pass tainted data to a sink node. Sanitizer nodes cleanse or untaint or endorse information to mediate across different levels of trust. Regular nodes do not return tainted data, and it is not an error to pass tainted data to regular nodes. If tainted data is passed to regular nodes, they merely propagate it to their successors without any mediation.

---

<sup>1</sup>We do not focus on *implicit* information flows [26] in this paper: discussions with CAT.NET [1] developers reveal that detecting explicit information flow vulnerabilities is a more urgent concern. Existing commercial tools in this space exclusively focus on explicit information flow.

```

1. void ProcessRequest(HttpRequest request,
2.     HttpResponse response)
3. {
4.     string s1 = request.GetParameter("name");
5.     string s2 = request.GetHeader("encoding");
6.
7.     response.WriteLine("Parameter " + s1);
8.     response.WriteLine("Header " + s2);
9. }

```

Figure 1: Simple cross-site scripting example

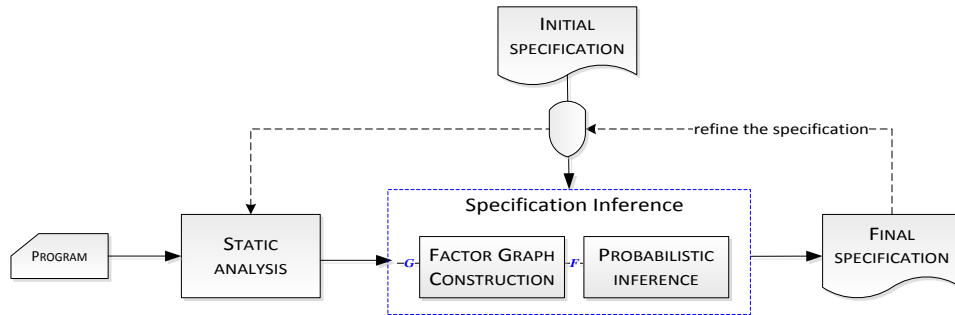


Figure 2: MERLIN system architecture.

A classification of nodes in a propagation graph into sources, sinks and sanitizers is called an *information flow specification* or, just *specification* for brevity. Given a propagation graph and a specification, one can easily run a reachability algorithm to check if all paths from sources to sinks pass through a sanitizer. In fact, this is precisely what many commercial analysis tools in everyday use do [5, 23].

User-provided specifications, however, lead to both false positives and false negatives in practice. False positives arise because a flow from source to sink classified as offending by the tool could have a sanitizer that the tool was unaware of. False negatives arise because of incomplete information about sources and sinks, as the following example illustrates.

This paper presents MERLIN, a tool that *automatically* infers information flow specifications for programs. Our inference algorithm uses the intuition that *most* paths in a propagation graph are secure. That is, most paths in the propagation graph that start from a source and end in a sink pass through some sanitizer.

**Example 1** Consider a Web application code snippet written in C# shown

in Figure 1. While method `GetParameter`, the method returning arguments of an HTTP request, is highly likely to be part of the default specification that comes with a static analysis tool and classified as a source, the method retrieving an HTTP header `GetHeader` may easily be missed. Because `response.WriteLine` sends information to the browser, there are two possibilities of *cross-site scripting vulnerabilities* on line 7 and line 8. The vulnerability in line 7 will be reported, but the vulnerability on line 8 may be missed due to an incomplete specification. In fact, in both .NET and J2EE there exist a number of source methods that return various parts of an HTTP request. When we run MERLIN on larger bodies of code, even within the `HttpRequest` class alone, MERLIN correctly determines that `getQueryString`, `getMethod`, `getEncoding`, and others are sources missing from the default specification that already contains 111 elements. While this example is small, it is meant to convey the substantial challenge involved in identifying appropriate APIs for an arbitrary application.  $\square$

**Our approach.** MERLIN infers information flow specifications using probabilistic inference. By using a random variable for each node in the propagation graph to denote whether the node is a source, sink or sanitizer, the intuition that most paths in a propagation graph are secure can be modeled using one *probabilistic constraint* for each path in the propagation graph. A probabilistic constraint is a path constraint parameterized by the probability that the constraint is true. By solving these constraints, we can get solutions to values of these random variables, which yields an information flow specification. In other words, we use probabilistic reasoning and the intuition we have about the outcome of the constraints (i.e, the probability of each constraint being true) to calculate values for the inputs to the constraints. Since there can be an exponential number of paths, using one constraint per path does not scale. In order to scale, we approximate the constraints using a different set of constraints on chosen triples of nodes in the graph. We show that the constraints on triples are a probabilistic abstraction of the constraints on paths (see Section 5) according to the theory developed by McIver and Morgan [20, 21]. As a consequence, we can show that approximation using constraints on triples does not introduce false positives when compared with the constraints on paths. After studying large applications, we found that we need additional constraints to reduce false positives, such as constraints to minimize the number of inferred sanitizers, and constraints to avoid inferring wrappers as sources or sinks. Section 2 describes these observations and constraints in detail. We show how to model these observations as additional probabilistic constraints. Once we have modeled the

problem as a set of probabilistic constraints, specification inference reduces to probabilistic inference. To perform probabilistic inference in a scalable manner, MERLIN uses *factor graphs*, which have been used in a variety of applications [13, 35]. While we can use the above approach to infer specifications *without* any prior specification, we find that the quality of inference is significantly higher if we use the default specification that comes with the static analysis tool as the initial specification, using MERLIN to “complete” this partial specification. Our empirical results in Section 6 demonstrate that our tool provides significant value in both these situations.

In our experiments, we use CAT.NET [1], a state-of-the-art static analysis tool for finding Web application security flaws that works on .NET bytecode. The initial specification is also modeled as extra probabilistic constraints on the random variables associated with nodes of the propagation graph for which we have partial specifications. To demonstrate the efficacy of MERLIN, we show empirical results from 10 large web applications where we have used MERLIN to complete existing partial specifications.

**Contributions.** Our paper makes the following contributions:

- MERLIN is the first practical approach to inferring specifications for explicit information flow analysis tools, a problem made important in recent years by the proliferation of information flow vulnerabilities in Web applications.
- A salient feature of our method is that our approximation (using triples instead of paths) can be characterized formally—we make a connection between probabilistic constraints and probabilistic programs, and use the theory of probabilistic refinement developed by McIver and Morgan [20, 21] to show refinement relationships between sets of probabilistic constraints. As a result, our approximation does not introduce false positives.
- MERLIN is able to successfully and efficiently infer information flow specifications in large code bases. We provide a comprehensive evaluation of the efficacy and practicality of MERLIN using 10 Web application benchmarks. We find a total of 167 new confirmed specifications, which result in a total of 302 vulnerabilities across the 10 benchmarks that have thus far remained undetected. MERLIN-inferred specifications also result in 13 false positives being removed.

**Outline.** The rest of the paper is organized as follows. Section 2 gives motivation for our specification inference techniques MERLIN uses and Section 3 gives background on factor graphs. Section 4 describes our algorithm

in detail. Section 5 proves that the system of triple constraints is a probabilistic abstraction over the system of path constraints. Section 6 provides our experimental evaluation. Finally, Sections 7 and 8 describe related work and conclude.

## 2 Overview

Figure 2 shows an architectural diagram of MERLIN. MERLIN starts with an initial, potentially incomplete specification of the application to produce a more complete specification. Returning to Example 1, MERLIN might start with an initial specification that classifies `GetParameter` as a source and `WriteLine` as a sink. The specification output by MERLIN additionally contains `GetHeader` as a source. In addition to the initial specification, MERLIN also consumes a *propagation graph*, a representation of the interprocedural data flow that is output by CAT.NET. MERLIN performs the following steps: 1) construct a factor graph based on the propagation graph; 2) Perform probabilistic inference on the factor graph to derive the likely specification.

The crux of our approach is probabilistic inference: we first use the propagation graph to generate a set of probabilistic constraints and then use factor graphs to efficiently solve them. We will use the rest of this section to describe the constraints and to build up the motivation for the rest of the paper.

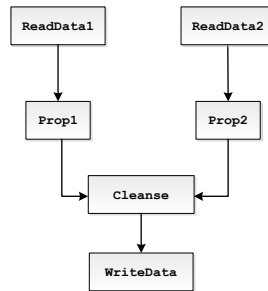
**Example 2** We shall illustrate the ideas in this section using the example program shown below.

```

public void TestMethod1() {
    string a = ReadData1();
    string b = Prop1(a);
    string c = Cleanse(b);
    WriteData(c);
}

public void TestMethod2() {
    string d = ReadData2();
    string e = Prop2(d);
    string f = Cleanse(e);
    WriteData(f);
}

```



In addition to two top-level “driver” methods, `TestMethod1` and `TestMethod2`, this code uses six methods: `ReadData1`, `ReadData2`, `Prop1`,



`Prop2`, `Cleanse`, and `WriteData`. This program gives rise to the propagation graph shown on the right. Intuitively, every node in this graph corresponds to a method. An edge  $m_1 \rightarrow m_2$  indicates that the return result of method  $m_1$  flows to method  $m_2$  as an argument.  $\square$

## 2.1 Assumptions and Probabilistic Inference

MERLIN relies on the assumption that *most* paths in the propagation graph are secure. That is, most paths that start from a source to a sink pass through a sanitizer node. The assumption that errors are rare has been used before in other specification inference techniques [4, 11]. Further, we assume that the number of sanitizers is small, relative to the number of regular nodes. Indeed, developers typically define a small number of sanitization functions or use ones supplied in libraries, and call them extensively. For instance, the out-of-box specification that comes with CAT.NET summarized in Figure 20 contains only 7 sanitizers.

However, as we show later in this section, applying these assumptions along various paths individually can lead to inconsistencies, since the constraints inferred from different paths can be mutually contradictory. Thus, we need to represent and analyze each path within a constraint system that tolerates uncertainty and contradictions. Therefore, we parameterize each constraint with the probability of its satisfaction. These probabilistic constraints model the relative positions of sources, sinks, and sanitizers in the propagation graph. Our goal is to classify each node as a source, sink, sanitizer, or a regular node, so as to optimize satisfaction of these probabilistic constraints. We solve these constraints using probabilistic inference.

## 2.2 Core Constraints

Figure 3 summarizes the constraints that MERLIN uses. We describe each of them in turn below referring to Example 2 where appropriate.

**B1: Path safety.** We believe that most paths from a source to a sink pass through at least one sanitizer. For example, we believe that if `ReadData1` is a source, and `WriteData` is a sink, at least one of `Prop1` or `Cleanse` is a sanitizer. This is stated using the set of constraints **B1** shown in Figure 3. While constraints **B1** models our core belief accurately, they are inefficient to use directly. **B1** requires one constraint per path, and the number of acyclic paths could be exponential in the number of nodes of the propagation graph.

**C1: Triple safety.** In order to abstract the constraint set **B1** with a polynomial number of constraints, we add a safety constraint **C1** for each

triple of nodes as shown in Figure 3. The number of **C1** constraints is  $O(N^3)$ , where  $N$  is the number of nodes in the propagation graph. In Section 5 we prove that the constraints **C1** are a probabilistic abstraction of constraints **B1** under suitable choices of parameters.

### 2.3 Auxiliary Constraints

In practice, the set of constraints **C1** does not limit the solution space enough. We have found empirically that just using this set of constraints allows too many possibilities, several of which are incorrect classifications. By looking at results over several large benchmarks we have come up with four sets of auxiliary constraints **C2**, **C3**, **C4**, and **C5**, which greatly enhance the precision of MERLIN inference.

**C2: Pairwise Minimization.** The set of constraints **C1** allows the solver

- 
- B1** For every acyclic path  $m_1, m_2, \dots, m_{k-1}, m_k$ , where  $m_1$  is a potential source and  $m_k$  is a potential sink, the joint probability of classifying  $m_1$  as a source,  $m_k$  as a sink and all of  $m_2, \dots, m_{k-1}$  as regular nodes is *low*.
  - C1** For every triple of nodes  $\langle m_1, m_2, m_3 \rangle$ , where  $m_1$  is a potential source,  $m_3$  is a potential sink, and  $m_1$  and  $m_3$  are connected by a path through  $m_2$  in the propagation graph, the joint probability that  $m_1$  is a source,  $m_2$  is not a sanitizer, and  $m_3$  is a sink is *low*.
  - C2** For every pair of nodes  $\langle m_1, m_2 \rangle$  such that both  $m_1$  and  $m_2$  lie on the same path from a potential source to a potential sink, the probability of both  $m_1$  and  $m_2$  being sanitizers is *low*.
  - C3** Each node  $m$  is classified as a sanitizer with probability  $s(m)$  (see Definition 2 for definition of  $s$ ).
  - C4** For every pair of nodes  $\langle m_1, m_2 \rangle$  such that both  $m_1$  and  $m_2$  are potential sources, and there is a path from  $m_1$  to  $m_2$  the probability that  $m_1$  is a source and  $m_2$  is not a source is *high*.
  - C5** For every pair of nodes  $\langle m_1, m_2 \rangle$  such that both  $m_1$  and  $m_2$  are potential sinks, and there is a path from  $m_1$  to  $m_2$  the probability that  $m_2$  is a sink and  $m_1$  is not a sink is *high*.
- 

Figure 3: Constraint formulation. Probabilities in italics are parameters of the constraints.

flexibility to assign multiple sanitizers along a path. In general, we want to minimize the number of sanitizers we infer. Thus, if there are several solutions to the set of constraints **C1**, we want to favor solutions that infer fewer sanitizers, while satisfying **C1**, with higher probability.

For instance, consider the path `ReadData1, Prop1, Cleanse, WriteData` in Example 2. Suppose `ReadData1` is a source and `WriteData` is a sink. **C1** constrains the triple

$$\langle \text{ReadData1}, \text{Prop1}, \text{WriteData} \rangle$$

so that the probability of `Prop1` not being a sanitizer is low; **C1** also constrains the triple

$$\langle \text{ReadData1}, \text{Cleanse}, \text{WriteData} \rangle$$

such that the probability of `Cleanse` not being a sanitizer is low. One solution to these constraints is to infer that both `Prop1` and `Cleanse` are sanitizers. In reality, programmers do not add multiple sanitizers on a path and we believe that only one of `Prop1` or `Cleanse` is a sanitizer. Thus, we add a constraint **C2** that for each pair of potential sanitizers it is unlikely that both are sanitizers, as shown in Figure 3. The number of **C2** constraints is  $O(N^2)$ , where  $N$  is the number of nodes in the propagation graph.

**Need for probabilistic constraints.** Note that constraints **C1** and **C2** can be mutually contradictory, if they are modeled as non-probabilistic boolean constraints. For example, consider the propagation graph of Example 2. With each of the nodes `ReadData1, WriteData, Prop1, Cleanse` let us associate boolean variables  $r_1, w, p_1$  and  $c$  respectively. The interpretation is that  $r_1$  is true iff `ReadData1` is source,  $w$  is true iff `WriteData` is a sink,  $p_1$  is true iff `Prop1` is a sanitizer, and  $c$  is true iff `Cleanse` is a sanitizer. Then, constraint **C1** for the triple  $\langle \text{ReadData1}, \text{Prop1}, \text{WriteData} \rangle$  is given by the boolean formula  $r_1 \wedge w \implies p_1$ , and the constraint **C1** for the triple  $\langle \text{ReadData1}, \text{Cleanse}, \text{WriteData} \rangle$  is given by the formula  $r_1 \wedge w \implies c$ . Constraint **C2** for the pair  $\langle \text{Prop1}, \text{Cleanse} \rangle$  states that both `Prop1` and `Cleanse` cannot be sanitizers, and is given by the formula  $\neg(p_1 \wedge c)$ . In addition, suppose we have additional information (say, from a partial specification given by the user) that `ReadData1` is indeed a source, and `WriteData` is a sink. We can conjoin all the above constraints to get the boolean formula:

$$(r_1 \wedge w \implies p_1) \wedge (r_1 \wedge w \implies c) \wedge \neg(p_1 \wedge c) \wedge r_1 \wedge w$$

This formula is unsatisfiable and these constraints are mutually contradictory. Viewing them as probabilistic constraints gives us the flexibility to

add such conflicting constraints, and the probabilistic inference algorithms resolve such conflicts by favoring satisfaction of those constraints with higher probabilities attached to them.

**C3: Sanitizer Prioritization.**

**Definition 1.** For each node  $m$  define weight  $W(m)$  to be the total number of paths from sources to sinks that pass through  $m$ .

Suppose we know that `ReadData1` is a source, `ReadData2` is *not* a source, and `WriteData` is a sink. Then  $W(\text{Prop1}) = W(\text{Cleanse}) = 1$ , since there is only one source-to-sink path that goes through each of them. However, in this case, we believe that `Prop1` is more likely to be a sanitizer than `Cleanse` since *all* paths going through `Prop1` are source-to-sink paths and only *some* paths going through `Cleanse` are source-to-sink paths.

**Definition 2.** For each node  $m$  define  $W_{total}(m)$  to be the total number of paths in the propagation graph that pass through the node  $m$  (this includes both source-to-sink paths, as well as other paths). Let us define  $s(m)$  for each node  $m$  as follows:

$$s(m) = \frac{W(m)}{W_{total}(m)}$$

We add a constraint **C3** that prioritizes each potential sanitizer based on its  $s()$  value, as shown in Figure 3. The number of **C3** constraints is  $O(N)$ , where  $N$  is the number of nodes in the propagation graph.

**C4: Source Wrapper Avoidance.** Similar to avoiding inference of multiple sanitizers on a path, we also wish to avoid inferring multiple sources on a path. A prominent issue with inferring sources is the issue of having wrappers, i.e. functions that return the result produced by the source. For instance, if an application defines their own series of wrappers around system APIs, which is not uncommon, there is no need to flag those as sources.

In such cases, we want MERLIN to infer the actual source rather than the wrapper function around it. We add a constraint **C4** for each pair of potential sources as shown in Figure 3. The number of **C4** constraints is  $O(N^2)$ , where  $N$  is the number of nodes in the propagation graph.

**C5: Sink Wrapper Avoidance.** Wrappers on sinks can be handled similarly, with the variation that in the case of sinks the data actually flows from the wrapper to the sink. We add a constraint **C5** for each pair of potential sinks as shown in Figure 3. The number of **C5** constraints is  $O(N^2)$ , where  $N$  is the number of nodes in the propagation graph.

### 3 Factor Graph Primer

In the previous section, we described a set of probabilistic boolean constraints that are generated from an input propagation graph. The conjunction of all these constraints can be looked upon as a joint probability distribution over random variables that measure the odds of propagation graph nodes being sources, sanitizers or sinks.

Let  $p(x_1, \dots, x_N)$  be a joint probability distribution over boolean variables  $x_1, \dots, x_N$ . We are interested in computing the *marginal probabilities*  $p_i(x_i)$  defined as:

$$p_i(x_i) = \sum_{x_1} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_N} p(x_1, \dots, x_N) \quad (1)$$

where  $x_i \in \{true, false\}$  for  $1 \leq i \leq N$ . Since there are an exponential number of terms in Equation 1, a naïve algorithm for computing  $p_i(x_i)$  will not work in practice. Equation 1 can also be written as

$$p_i(x_i) = \sum_{\sim\{x_i\}} p(x_1, \dots, x_N) \quad (2)$$

where the sum is over all variables except  $x_i$ . The marginal probability for each variable defines the solution that we are interested in computing. Intuitively, these marginals correspond to the likelihood of each boolean variable being equal to *true* or *false*.

Factor graphs [35] are graphical models that are used for computing marginal probabilities efficiently. These graphs take advantage of their structure in order to speed up the marginal probability computation (known as probabilistic inference). There are a wide variety of techniques for performing probabilistic inference on a factor graph and the *sum-product* algorithm [35] is the most practical algorithm among these.

Let the joint probability distribution  $p(x_1, \dots, x_N)$  be a product of factors as follows:

$$p(x_1, \dots, x_N) = \prod_s f_s(x_s) \quad (3)$$

where  $x_s$  is the set of variables involved in the factor  $f_s$ . A *factor graph* is a bipartite graph that represents this factorization. A factor graph has two types of nodes:

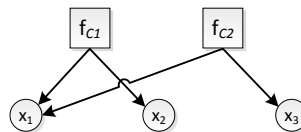


Figure 4: Factor graph for (3).

- *Variable nodes*: one node for every variable  $x_i$ .
- *Function nodes*: one node for every function  $f_s$ .

**Example 3** As an example, consider the following formula

$$\underbrace{(x_1 \vee x_2)}_{C_1} \wedge \underbrace{(\neg x_1 \vee x_3)}_{C_2} \quad (4)$$

Equation 4 can be rewritten as:

$$f(x_1, x_2, x_3) = f_{C_1}(x_1, x_2) \wedge f_{C_2}(x_1, x_3) \quad (5)$$

where

$$f_{C_1} = \begin{cases} 1 & \text{if } x_1 \vee x_2 = \text{true} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

$$f_{C_2} = \begin{cases} 1 & \text{if } x_1 \vee \neg x_3 = \text{true} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

The factor graph for this formula is shown in Fig. 4. There are three variable nodes for each variable  $x_i$ ,  $1 \leq i \leq 3$  and a function node  $f_{C_j}$  for each clause  $C_j$ ,  $j \in \{1, 2\}$ .

Instead of using boolean values to check satisfiability of formula 4, we could use probabilistic inference: then we must interpret the formula as a probabilistic constraint.

$$p(x_1, x_2, x_3) = \frac{f_{C_1}(x_1, x_2) \times f_{C_2}(x_1, x_3)}{Z}, \quad (8)$$

where

$$Z = \sum_{x_1, x_2, x_3} (f_{C_1}(x_1, x_2) \times f_{C_2}(x_1, x_3)) \quad (9)$$

is the normalization constant. The marginal probabilities are defined as

$$p_i(x_i) = \sum_{\sim\{x_i\}} p(\vec{x}), \quad 1 \leq i \leq 3 \quad (10)$$

Equations 6 and 7 can also be defined probabilistically thus allowing for solutions that do not satisfy formula 4; but such solutions are usually set up such that they occur with low probability as shown below:

$$f_{C_1} = \begin{cases} 0.9 & \text{if } x_1 \vee x_2 = \text{true} \\ 0.1 & \text{otherwise} \end{cases} \quad (11)$$

$$f_{C_2} = \begin{cases} 0.9 & \text{if } x_1 \vee \neg x_3 = \text{true} \\ 0.1 & \text{otherwise} \end{cases} \quad (12)$$

$p_i(x_i = \text{true})$  denotes the fraction of solutions where the variable  $x_i$  has value *true*.  $\square$

This information can be used to find a solution efficiently – choose a variable with the highest probability and assign its value as specified by the marginal probability and recompute marginal probabilities and iterate until all variables have been assigned values. Exact details of iterative approaches that operate on factor graphs in order to compute marginal probabilities can be found in [13, 35].

## 4 Inference Algorithm

As described in Section 2 and shown in Figure 2, MERLIN’s processing progresses in stages. The first step in MERLIN is to construct a propagation graph.

**Definition 3.** *A propagation graph is a directed graph  $G = \langle N_G, E_G \rangle$ , where nodes  $N_G$  are methods and an edge  $(n_1 \rightarrow n_2) \in E_G$  indicates that there is a flow of information from method  $n_1$  to method  $n_2$  through method arguments, or return values, or indirectly through pointers.*

The propagation graph is a representation of the interprocedural data flow produced by CAT.NET, the static analysis tool for .NET. Further steps of MERLIN use the propagation graph as a representation of information flow between methods in the program.

Section 4.2 describes how we convert a propagation graph  $P$  into a factor graph  $F$  and in Section 4.3, we describe the computation of the functions  $W()$  and  $s()$  (from Definitions 1 and 2 in Section 2) that are used in the construction of the factor graph.

### 4.1 Reducing the Propagation Graph

The reduced propagation graph abstracts away the intraprocedural details and creates a considerably smaller interprocedural data flow representation to operate on. Figure 5 shows the original propagation graph for Example 2. Because of space limitations, the figure only shows function `TestMethod1`. The edges within the thick dashed boundary are intraprocedural edges; the rest of the edges crossing the boundary are interprocedural.

As shown in Figure 7, to construct a compressed propagation graph, we systematically traverse dataflow paths in graph  $P$  to detect when the return result of one method is passed to another.

## 4.2 Constructing the Factor Graph

Figure 9 shows the algorithm `GenFactorGraph` that we use to generate the factor graph from the propagation graph. The construction of the factor graph proceeds as follows. First, in line 1, the procedure `MakeAcyclic` converts the input propagation graph into a DAG  $G'$ , by doing a breadth first search, and deleting edges that close cycles. Next, in line 2, the procedure `ComputePairsAndTriples` computes four sets shown in Figure 8.

These sets can be computed by doing a topological sort of  $G'$  (the acyclic graph), and doing one pass over the graph in topological order and recording for each internal node, the set of sources and sinks that can be reached from the node. These sets can be computed in  $O(N^3)$  time, where  $N$  is the number of nodes in the propagation graph.

Next in line 3, the function `ComputeWAndSValues` is invoked to compute the  $W(n)$  and  $s(n)$  for every potential sanitizer  $n$ . The function `ComputeWAndSValues` is described in Figure 10. In lines 4–8, the algorithm creates a factor node for the constraints **C1**. In lines 9–13, the algorithm iterates through all pairs  $\langle b_1, b_2 \rangle$  of potential sanitizers (that is, actual sanitizers as well as regular nodes) such that there is a path in the propagation graph from  $b_1$  to  $b_2$  and adds factors for constraints **C2**. In lines 14–18, the algorithm iterates through all potential sanitizers and adds factors for constraints **C3**. In lines 19–23, the algorithm iterates through all pairs  $\langle a_1, a_2 \rangle$  of potential sources such that there is a path in the propagation graph from  $a_1$  to  $a_2$  and adds factors for constraints **C4**. Similarly, in lines 24–28, the algorithm iterates through all pairs  $\langle c_1, c_2 \rangle$  of potential sinks such that there is a path in the propagation graph from  $c_1$  to  $c_2$  and adds factors for constraints **C5**.

**Example 4** Figure 6 shows the factor graph obtained by applying algorithm `FactorGraph` to the propagation graph in Figure 2. The marginal probabilities for all variable nodes are computed by performing probabilistic inference on the factor graph and these are used to classify sources, sanitizers and sinks in the propagation graph.  $\square$



### 4.3 Computing $s()$ and $W()$

Recall the definitions of  $s()$  and  $W()$  from Definitions 1 and 2. Figure 10 shows the body of `ComputeWAndSValues`, which computes  $s(n)$  for each potential sanitizer node, given input probabilities for each potential source and each potential sink. The value  $s(n)$  for each potential sanitizer  $n$  depends on the sum of weighted source-sink paths that go through  $n$  divided by the total number of paths that go through  $n$ . The algorithm computes  $W(n)$  and  $s(n)$  by computing four numbers  $F(n)$ ,  $F_{total}(n)$ ,  $B(n)$  and  $B_{total}(n)$ .

$F(n)$  denotes the total number of sources that can reach  $n$ , and  $F_{total}(n)$  denotes the total number of paths that can reach  $n$ .  $B(n)$  denotes the total number of sinks that can be reached from  $n$  and  $B_{total}(n)$  denotes the total number of paths that can be reached from  $n$ . Since the graph  $G'$  is a DAG,  $F(n)$  and  $F_{total}(n)$  can be computed by traversing potential sanitizers in topological sorted order, and  $B(n)$  and  $B_{total}(n)$  can be computed by traversing potential sanitizers in reverse topological order. The computation of  $F(n)$  and  $F_{total}(n)$  in forward topological order is done in lines 5–12 and the computation of  $B(n)$  and  $B_{total}(n)$  in reverse topological order is done in lines 17–24. Once  $F(n)$  and  $B(n)$  are computed,  $W(n)$  is set to  $F(n) \times B(n)$  and  $s(n)$  is set to

$$\frac{W(n)}{F_{total}(n) \times B_{total}(n)} = \frac{W(n)}{W_{total}(n)}$$

as shown in line 26.

## 5 Relationship between Triples and Paths

In this section, we give a formal relationship between the exponential number of constraints **B1** and the cubic number of constraints **C1** in Section 2. In particular, we use the theory of probabilistic abstraction and refinement developed by McIver and Morgan [20, 21] to derive appropriate bounds on probabilities associated with constraints **B1** and **C1** so that **C1** is a probabilistic abstraction of the specification **B1** (or, equivalently, **B1** is a probabilistic refinement of **C1**). We first introduce some terminology and basic concepts from [21].

**Probabilistic refinement primer.** Non-probabilistic programs can be reasoned with assertions in the style of Floyd and Hoare [8]. The following formula in Hoare logic:

$$\{Pre\} \text{ Prog } \{Post\}$$

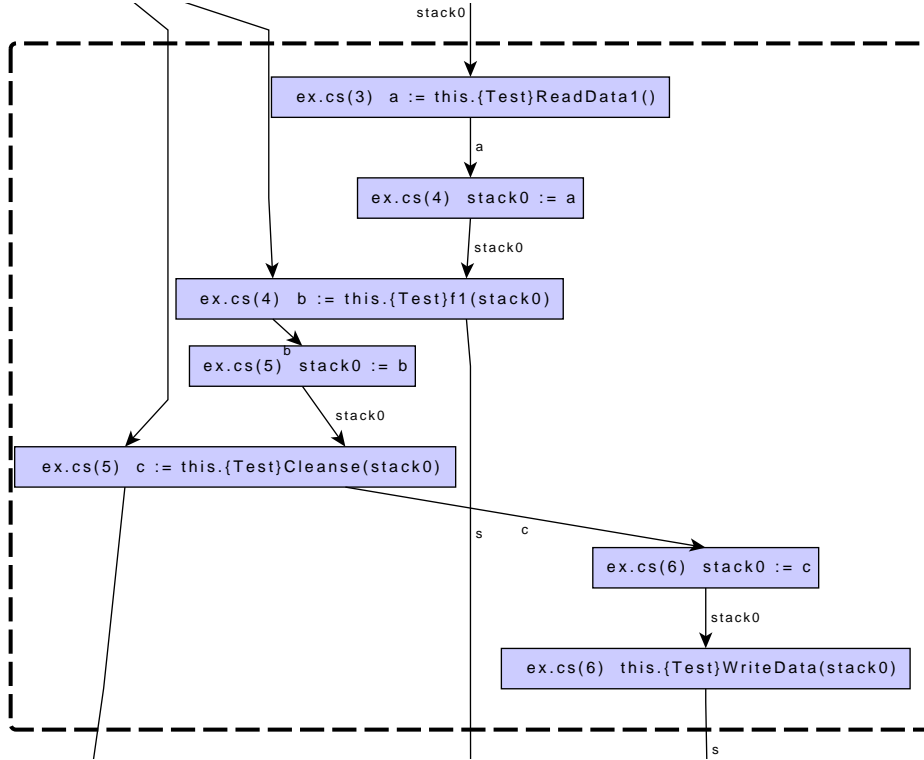


Figure 5: Original propagation graph for Example 2.

is valid if for every state  $\sigma$  satisfying the assertion  $Pre$ , if the program  $Prog$  is started at  $\sigma$ , then the resulting state  $\sigma'$  satisfies the assertion  $Post$ . We assume  $Prog$  always terminates, and thus we do not distinguish partial and total correctness.

McIver and Morgan extend such reasoning to probabilistic programs [20, 21]. In order to reason about probabilistic programs, they generalize asser-

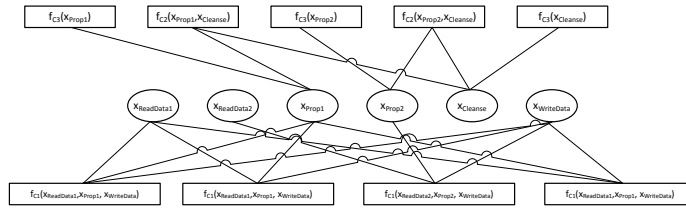


Figure 6: Factor graph for the propagation graph in Example 2.

tions to expectations. An *expectation* is a function that maps each state to a positive real number. If **Prog** is a probabilistic program, and  $Pre_E$  and  $Post_E$  are expectations, then the probabilistic Hoare-triple

$$\{Pre_E\} \text{ Prog } \{Post_E\}$$

is interpreted to mean the following: If the program **Prog** is started with an initial expectation  $Pre_E$ , then it results in the expectation  $Post_E$  after execution.

Assertions are ordered by implication ordering. Expectations are ordered by the partial order  $\Rightarrow$ . Given two expectations  $A_E$  and  $B_E$ , we say that  $A_E \Rightarrow B_E$  holds if for all states  $\sigma$ , we have that  $A_E(\sigma) \leq B_E(\sigma)$ . Given an assertion  $A$  the expectation  $[A]$  is defined to map every state  $\sigma$  to 1 if  $\sigma$  satisfies  $A$  and to 0 otherwise.

BuildReducedPropGraph

Inputs:

Propagation graph  $P$

Returns:

Reduced propagation graph  $G$

```

1: for all node  $n \in P$  do
2:   if  $n.IsCall$  then
3:     for all  $m \in n.Callees$  do
4:       TRAVERSE( $m, n, \emptyset, true$ )
5:     end for
6:   end if
7: end for

{This subroutine is called repeatedly}
8: function TRAVERSE( $s, n, visited, allowToExit$ )
9: if  $n \notin visited$  then
10:   $visited = visited \cup \{n\}$ 
11:  for all  $e \in n.OutgoingEdges$  do
12:    if  $\neg allowToExit \wedge e.Type.IsInterprocedural$  then
13:      continue
14:    end if
15:    if  $e.Type.IsInterprocedural$  then
16:      if  $e.To.IsCall$  then
17:        for  $t \in e.To.Callees$  do
18:           $G.Edges = G.Edges \cup \{s \rightarrow t\}$ 
19:        end for
20:      else
21:        TRAVERSE( $s, e.To, visited, e.Type.IsInterprocedural$ );
22:      end if
23:    end if
24:  end for
25: end if

```

Figure 7: Constructing reduced propagation graph  $G$  from the original propagation graph  $P$ .

Set	Definition
$Triples$	$\bigcup_{p \in paths(G')} \{ \langle x_{src}, x_{san}, x_{snk} \rangle \mid \begin{array}{l} x_{src} \text{ is connected to} \\ x_{snk} \text{ via } x_{san} \text{ in } p \end{array} \}$
$Pairs_{src}$	$\bigcup_{p \in paths(G')} \{ \langle x_{src}, x'_{src} \rangle \mid \begin{array}{l} x_{src} \text{ is connected to} \\ x'_{src} \text{ in } p \end{array} \}$
$Pairs_{san}$	$\bigcup_{p \in paths(G')} \{ \langle x_{san}, x'_{san} \rangle \mid \begin{array}{l} x_{san} \text{ is connected to} \\ x'_{san} \text{ in } p \end{array} \}$
$Pairs_{snk}$	$\bigcup_{p \in paths(G')} \{ \langle x_{snk}, x'_{snk} \rangle \mid \begin{array}{l} x_{snk} \text{ is connected to} \\ x'_{snk} \text{ in } p \end{array} \}$

Figure 8: Set definitions for algorithm in Figure 9.

Suppose  $A_E \Rightarrow B_E$ . Consider a sampler that samples states using the expectations as a probability measure. Then, for any threshold  $t$  and state  $\sigma$ , if  $A_E(\sigma) > t$ , then it is the case that  $B_E(\sigma) > t$ . In other words, for any sampler with any threshold  $t$ , sampling over  $A_E$  results in a subset of states than those obtained by sampling over  $B_E$ .

Traditional axiomatic proofs are done using weakest preconditions. The weakest precondition operator is denoted by  $WP$ . By definition, for any program  $\mathbf{Prog}$  and assertion  $A$ , we have that  $WP(\mathbf{Prog}, A)$  to be the weakest assertion  $B$  (weakest is defined with respect to the implication ordering between assertions) such that the Hoare triple  $\{B\}\mathbf{Prog}\{A\}$  holds.

McIver and Morgan extend weakest preconditions to expectations, and define for an expectation  $A_E$ , and a probabilistic program  $\mathbf{Prog}$ ,  $WP(\mathbf{Prog}, A_E)$  is the weakest expectation  $B_E$  (weakest is defined with respect to the ordering  $\Rightarrow$  between expectations) such that the probabilistic Hoare triple  $\{B_E\}\mathbf{Prog}\{A_E\}$  holds. Given two probabilistic programs  $\mathbf{Spec}$  and  $\mathbf{Impl}$  with respect to a post expectation  $Post_E$ , we say that  $\mathbf{Impl}$  *refines*  $\mathbf{Spec}$  if  $WP(\mathbf{Spec}, Post_E) \Rightarrow WP(\mathbf{Impl}, Post_E)$ .

**Refinement between constraint systems.** We now model constraints **B1** and **C1** from Section 2 as probabilistic programs with an appropriate post expectation, and derive relationships between the parameters of **B1** and **C1** such that **B1** refines **C1**.

Consider any directed acyclic graph  $G = \langle V, E \rangle$ , where  $E \subseteq V \times V$ . In this simple setting, nodes with indegree 0 are potential sources, nodes with outdegree 0 are potential sinks, and other nodes (internal nodes with both indegree and outdegree greater than 0) are potential sanitizers. We want

to classify every node in  $V$  with a boolean value 0 or 1. That is, we want a mapping  $m : V \rightarrow \{0, 1\}$ , with the interpretation that for a potential source  $s \in V$ ,  $m(s) = 1$  means that  $s$  is classified as a source, and that for a potential sink  $n \in V$ ,  $m(n) = 1$  means that  $n$  is classified as a sink, and that for a potential sanitizer  $w \in V$ ,  $m(w) = 1$  means that  $w$  is classified as a sanitizer. We extend the mapping  $m$  to operate on paths (triples) over  $G$  by applying  $m$  to every vertex along the path (triple).

We want mappings  $m$  that satisfy the constraint that for any path  $p = s, w_1, w_2, \dots, w_m, n$  that starts at a potential source  $s$  and ends in a potential sink, the string  $m(p) \notin 10^*1$ , where  $10^*1$  is the language of strings that begin and end with 1 and have a sequence of 0's of arbitrary length in between.

The constraint set **B1** from Section 2 is equivalent in this setting to

**GenFactorGraph**

**Inputs:**

$(G = \langle V, E \rangle : \text{PropagationGraph})$ ,

parameters  $low_1, low_2, high_1, high_2, high_3, high_4 \in [0..1]$

**Returns:**

a factor graph  $F$  for the propagation graph  $G$

```

1:  $G' = \text{MakeAcyclic}(G)$ 
2:  $\langle \text{Triples}, \text{Pairs}_{src}, \text{Pairs}_{san}, \text{Pairs}_{snk} \rangle = \text{ComputePairsAndTriples}(G')$ 
3:  $s = \text{ComputeWAndSValues}(G')$ 
4: for each triple  $\langle a, b, c \rangle \in \text{Triples}$  do
5:   Create a factor  $f_{C1}(x_a, x_b, x_c)$  in the factor graph
6:   Let  $f_{C1}(x_a, x_b, x_c) = x_a \wedge \neg x_b \wedge x_c$ 
7:   Let probability  $\Pr(f_{C1}(x_a, x_b, x_c) = \text{true}) = low_1$ 
8: end for
9: for each pair  $\langle b_1, b_2 \rangle \in \text{Pairs}_{san}$  do
10:  Create a factor  $f_{C2}(x_{b_1}, x_{b_2})$  in the factor graph
11:  Let  $f_{C2}(x_{b_1}, x_{b_2}) = x_{b_1} \wedge x_{b_2}$ 
12:  Let probability  $\Pr(f_{C2}(x_{b_1}, x_{b_2}) = \text{true}) = low_2$ 
13: end for
14: for each potential sanitizer  $n \in V$  do
15:  Create a factor  $f_{C3}(x_n)$  in the factor graph
16:  Let  $f_{C3}(x_n) = x_n$ 
17:  Let  $\Pr(f_{C3}(x_n) = \text{true}) = s(n)$ 
18: end for
19: for each pair  $\langle x_{a_1}, x_{a_2} \rangle \in \text{Pairs}_{src}$  do
20:  Create a factor  $f_{C4}(x_{a_1}, x_{a_2})$  in the factor graph
21:  Let  $f_{C4}(x_{a_1}, x_{a_2}) = x_{a_1} \wedge \neg x_{a_2}$ 
22:  Let probability  $\Pr(f_{C4}(x_{a_1}, x_{a_2}) = \text{true}) = high_3$ 
23: end for
24: for each pair  $\langle x_{c_1}, x_{c_2} \rangle \in \text{Pairs}_{snk}$  do
25:  Create a factor  $f_{C5}(x_{c_1}, x_{c_2})$  in the factor graph
26:  Let  $f_{C5}(x_{c_1}, x_{c_2}) = \neg x_{c_1} \wedge x_{c_2}$ 
27:  Let probability  $\Pr(f_{C5}(x_{c_1}, x_{c_2}) = \text{true}) = high_4$ 
28: end for

```

Figure 9: Generating a factor graph from a propagation graph.

the probabilistic program `Path` given in Figure 16, and the constraint set **C1** from Section 2 is equivalent in this setting to the probabilistic program `Triple` given in Figure 17. The statement `assume(e)` is a no-op if  $e$  holds and silently stops execution if  $e$  does not hold. The probabilistic statement  $S_1 \oplus_q S_2$  executes statement  $S_1$  with probability  $q$  and statement  $S_2$  with probability  $1 - q$ . Note that both programs `Path` and `Triple` have the same post expectation  $[\forall \text{ paths } p \text{ in } G, m(p) \notin 10^*1]$ . Further, note that both programs are parameterized. The `Path` program has a parameter  $c_p$  associated with each path  $p$  in  $G$ , and the `Triple` program has a parameter  $c_t$  associated with each triple  $t$  in  $G$ .

The following theorem states that the probabilistic program `Path` re-

```

ComputeWAndSValues( $G : PropagationGraph$ )
Precondition:
Inputs:
acyclic propagation graph  $G$ 
Returns:
 $W(n)$  and  $s(n)$  for each potential sanitizer  $n$  in  $G$ 

1: for each potential source  $n \in V$  do
2:    $F(n) :=$  probability of  $n$  being a source node
3:    $F_{total}(n) := 1$ 
4: end for
5: for each potential sanitizer  $n \in V$  in topological order do
6:    $F(n) := 0$ 
7:    $F_{total}(n) := 0$ 
8:   for each  $m \in V$  such that  $(m, n) \in E$  do
9:      $F(n) := F(n) + F(m)$ 
10:     $F_{total}(n) := F_{total}(n) + F_{total}(m)$ 
11:   end for
12: end for
13: for each potential sink  $n \in V$  do
14:    $B(n) :=$  probability of  $n$  being a sink node
15:    $B_{total}(n) := 1$ 
16: end for
17: for each potential sanitizer  $n \in V$  in reverse topological order do
18:    $B(n) := 0$ 
19:    $B_{total}(n) := 0$ 
20:   for each  $m \in V$  such that  $(n, m) \in E$  do
21:      $B(n) := B(n) + B(m)$ 
22:      $B_{total}(n) := B_{total}(n) + B_{total}(m)$ 
23:   end for
24: end for
25: for each potential sanitizer  $n \in V$  do
26:    $W(n) := F(n) * B(n)$ 
27:    $s(n) := \frac{W(n)}{F_{total}(n) * B_{total}(n)}$ 
28: end for
29: return  $s$ 

```

Figure 10: Computing  $W(n)$  and  $s(n)$ .

Benchmark	LOC	DLLs	DLL sz
Alias Management Tool	10,812	3	65
Bicycle Club App	14,529	3	62
Software Catalog	11,941	15	118
Sporting Field Management Tool	15,803	3	290
Commitment Management Tool	25,602	7	369
New Hire Tool	5,595	11	565
Expense Report Approval Tool	78,914	4	421,200
Relationship Management	1,810,585	5	3,345
Customer Support Portal	66,385	14	2,447

Figure 11: Benchmark statistics.

finest program Triple under appropriate choices of probabilities as parameters. Furthermore, given a program Path with arbitrary values for the parameters  $c_p$  for each path  $p$ , it is possible to choose parameter values  $c_t$  for each triple  $t$  in the program Triple such that Path refines Triple.

**Theorem.** Consider any directed acyclic graph  $G = \langle V, E \rangle$  and probabilistic programs Path (Figure 16) and Triple (Figure 17) with stated post expectations. Let the program Path have a parameter  $c_p$  for each path  $p$ . For any such valuations to the  $c_p$ 's there exist parameter values for the Triple program, namely a parameter  $c_t$  for each triple  $t$  such that the program Path refines the program Triple with respect to the post expectation  $Post_E = [\forall \text{ paths } p \text{ in } G, m(p) \notin 10^*1]$ .

Benchmark	$G$		$F$	
	Nodes	Edges	Vars	Nodes
Alias Management Tool	59	1,209	3	3
Terralever	156	187	25	33
Bicycle Club App	176	246	70	317
Software Catalog	190	455	73	484
Sporting Field Management Tool	268	320	50	50
Commitment Management Tool	356	563	107	1,781
New Hire Tool	502	1,101	116	1,917
Expense Report Approval Tool	811	1,753	252	2,592
Relationship Management	3,639	22,188	874	391,221
Customer Support Portal	3,881	11,196	942	181,943

Figure 12: Size statistics for the propagation graph  $G$  and factor graph  $F$  used by MERLIN.

Benchmark	SOURCES				SANITIZERS				SINKS				
	All	✓	?	✗	All	✓	?	✗	All	✓	?	✗	Rate
Alias Management Tool	0	0	0	0	0	0	0	0	0	0	0	0	N/A
Terralever	1	1	0	0	0	0	0	0	0	2	0	0	0%
Bicycle Club App	11	11	0	0	3	2	0	1	7	4	0	3	33%
Software Catalog	1	1	0	0	8	3	0	5	6	3	2	1	62%
Sporting Field Management Tool	0	0	0	0	0	0	0	0	1	0	1	0	N/A
Commitment Management Tool	20	19	0	1	9	1	2	6	11	8	1	2	5%
New Hire Tool	3	3	0	0	1	1	0	0	17	14	0	3	0%
Expense Report Approval Tool	8	8	0	0	20	2	13	5	20	14	0	6	0%
Relationship Management	44	3	36	5	1	0	0	1	4	0	3	1	11%
Customer Support Portal	26	21	4	1	39	16	10	13	118	30	55	33	3%
<b>Total</b>	114	67	40	7	186	25	25	81	186	75	62	49	6%
													43%
													26%

Figure 13: New specifications discovered with MERLIN.



Benchmark	BEFORE				AFTER				
	All	✓	?	✗	All	✓	-	?	✗
Alias Management Tool	2	2	0	0	2	2	0	0	0
Terralever	0	0	0	0	1	1	0	0	0
Bicycle Club App	0	0	0	0	4	3	0	1	0
Software Catalog	14	8	0	6	8	8	6	0	0
Sporting Field Management	0	0	0	0	0	0	0	0	0
Commitment Management Tool	1	1	0	0	22	16	0	3	3
New Hire Tool	4	4	0	0	3	3	1	0	0
Expense Report Approval Tool	0	0	0	0	2	2	0	0	0
Relationship Management	9	6	3	0	13	10	3	0	0
Customer Support Portal	59	19	3	37	280	277	3	13	3
<b>Total</b>	89	40	6	43	335	342	13	17	3

Figure 14: Vulnerabilities before and after MERLIN.

*Proof:* Consider any triple  $t = \langle s, w, n \rangle$ . Choose the parameter  $c_t$  for the triple  $t$  to be equal to the product of the parameters  $c_p$  of all paths  $p$  in  $G$  that start at  $s$ , end at  $n$  and go through  $w$ . That is,

$$c_t = \prod_p c_p \quad (13)$$

such that  $t$  is a subsequence of  $p$ .

To show that Path refines Triple with respect to the post expectation  $Post_E$  stated in the theorem, we need to show that  $WP(\text{Triple}, Post_E) \Rightarrow$

Benchmark	CAT.NET	MERLIN		Total time
	$P$	$G$	$F$	
Alias Management Tool	2.64	4.59	2.63	9.86
Terralever	4.61	.81	2.67	8.09
Bicycle Club App	2.81	.53	2.72	6.06
Software Catalog	3.94	1.02	2.73	7.69
Sporting Field Management Tool	5.97	2.22	2.69	10.88
Commitment Management Tool	6.41	18.84	2.91	28.16
New Hire Tool	7.84	2.98	3.44	14.27
Expense Report Approval Tool	7.27	3.59	3.05	13.91
Relationship Management	55.38	87.63	66.45	209.45
Customer Support Portal	89.75	29.75	31.55	151.05

Figure 15: Static analysis and specification inference running time, in seconds.

$\text{Path}(G = \langle V, E \rangle)$   
 Returns:  
 Mapping  $m$  from  $V$  to the set  $\{0, 1\}$

1: for all paths  $p = s, \dots, n$  from potential sources to sinks in  $G$  do  
 2:    $\text{assume}(m(p) \notin 10^*1) \oplus_{c_p} \text{assume}(m(p) \in 10^*1)$   
 3: end for  
 Post expectation:  $[\forall \text{ paths } p \text{ in } G, m(p) \notin 10^*1]$ .

Figure 16: Algorithm Path

$\text{Triple}(G = \langle V, E \rangle)$   
 Returns:  
 Mapping  $m$  from  $V$  to the set  $\{0, 1\}$

1: for all triples  $t = \langle s, w, n \rangle$  such that  $s$  is a potential source,  $n$  is a potential sink and  $w$  lies on some path from  $s$  to  $n$  in  $G$  do  
 2:    $\text{assume}(m(\langle s, w, n \rangle) \neq 101) \oplus_{c_t} \text{assume}(m(\langle s, w, n \rangle) = 101)$   
 3: end for  
 Post expectation:  $[\forall \text{ paths } p \text{ in } G, m(p) \notin 10^*1]$ .

Figure 17: Algorithm Triple

$\text{WP}(\text{Path}, \text{Post}_E)$ . That is, for each state  $\sigma$ , we need to show that  $\text{WP}(\text{Triple}, \text{Post}_E)(\sigma) \leq \text{WP}(\text{Path}, \text{Post}_E)(\sigma)$ .

Note that  $\text{WP}(\text{assume}(e), [A]) = [e \wedge A]$ , and  $\text{WP}(S_1 \oplus_q S_2, [A]) = q \times \text{WP}(S_1, [A]) + (1 - q) \times \text{WP}(S_2, [A])$  [21]. Using these two rules, we can compute  $\text{WP}(\text{Triple}, \text{Post}_E)$  and  $\text{WP}(\text{Path}, \text{Post}_E)$  as an expression tree which is a sum of product of expressions, where each product corresponds to a combination of probabilistic choices made in the program.

First, consider any state  $\sigma$  that does not satisfy  $\text{Post}_E$ . For this state,  $\text{WP}(\text{Triple}, \text{Post}_E)(\sigma) = \text{WP}(\text{Path}, \text{Post}_E)(\sigma) = 0$ , and the theorem follows trivially. Next, consider a state  $\omega$  that satisfies  $\text{Post}_E$ . In this case,  $\text{WP}(\text{Path}, \text{Post}_E)(\omega)$  is the product of probabilities  $c_p$  for each path  $p$  in  $G$ .

Also, in this case  $\text{WP}(\text{Triple}, \text{Post}_E)(\omega)$  is the product of two quantities  $X(\omega)$  and  $Y(\omega)$ , where  $X(\omega)$  is equal to the product of probabilities  $c_t$  for each triple  $t = \langle s, w, n \rangle$  such that  $m(\langle s, w, n \rangle) \neq 101$ , and  $Y(\omega)$  is equal to the product of probabilities  $(1 - c_{t'})$  for each triple  $t' = \langle s', w', n' \rangle$  such that  $m(\langle s', w', n' \rangle) = 101$ . Since  $c_t$ 's have been carefully chosen according to Equation 13 and  $Y(\omega) \in [0, 1]$ , it follows that  $X(\omega)$  is less than or equal to the product of the probabilities  $c_p$  for each path  $p$ . Therefore, it is indeed the case that for each state  $\omega$ ,  $\text{WP}(\text{Triple}, \text{Post}_E)(\omega) \leq \text{WP}(\text{Path}, \text{Post}_E)(\omega)$ . ■

Any solver for a probabilistic constraint system  $C$  with post expectation

$Post_E$  chooses states  $\sigma$  such that  $WP(C, Post_E)(\sigma)$  is greater than some threshold  $t$ . Since we have proved that **Path** refines **Triple**, we know that every solution state for the **Triple** system, is also a solution state for the **Path** system. Thus, the set of states that are chosen by solver for the **Triple** system is contained in the set of states that are chosen by the solver for the **Path** system. *This has the desirable property that the **Triple** system will not introduce more false positives than the **Path** system.*

Note that the **Path** system itself can result in false positives, since it requires at least one sanitizer on each source-sink path, and does not require minimization of sanitizers. In order to remove false positives due to redundant sanitizers, we add the constraints **C2** and **C3** to the **Triple** system. Further, the path system does not distinguish wrappers of sources or sinks, so we add additional constraints **C4** and **C5** to avoid classifying these wrappers as sources or sinks. Using all these extra constraints, we find that the **Triple** system performs very well on several large benchmarks and infers specifications with very few false errors. We describe these results in the next section.

## 6 Experimental Evaluation

We believe that the ultimate success of a specification inference project such as MERLIN is to be judged on the quality of the experimental results it produces. This section presents the results of evaluating MERLIN on 10 large .NET Web applications. All of these are enterprise line of business applications currently in production written in C# on top of ASP.NET.

### 6.1 Experimental Setup

Figure 11 summarizes information about our benchmarks. Column 2 shows the traditional line-of-code metric for *all* the code within the application. However, as we discovered, not all code is actually deployed to the Web server, and the number and size of deployed DLLs primarily consisting of .NET bytecode is more relevant. These applications compile to one or more DLLs; the number of DLLs is shown in column 3. The total size of these DLLs, as measured in KB, is shown in column 4.

To put our results on specification discovery in perspective, Figure 20 provides information about the out-of-the box specification for CAT.NET, the static analysis tool that we used for our experiments. The second column shows the number of specifications for each specification type. The last column shows the number of *revisions* each portion of the specification has

```

1 public static string CreateQueryLink(
2     HttpRequest request, string key, string value,
3     List<string> keysToOmit, bool ignoreKey)
4 {
5     StringBuilder builder = new StringBuilder(
6         request.Url.AbsolutePath);
7     if (keysToOmit == null) {
8         keysToOmit = new List<string>();
9     }
10    builder.Append("?");
11    for (int i = 0; i < request.QueryString.Count; i++) {
12        if ((request.QueryString.GetKey(i) != key) &&
13            !keysToOmit.Contains(request.QueryString.GetKey(i)))
14        {
15            builder.Append(request.QueryString.GetKey(i));
16            builder.Append("=");
17            builder.Append(AntiXss.UrlEncode(
18                QueryStringParser.Parse(
19                    request.QueryString.GetKey(i))));
20            builder.Append("&");
21        }
22    }
23    if (!ignoreKey) {
24        builder.Append(key);
25        builder.Append("=");
26        builder.Append(AntiXss.UrlEncode(value));
27    }
28    return builder.ToString().TrimEnd(new char[] { '&' });
29 }

```

Figure 18: Function CreateQueryLink for Example 5.

---

```

Sources      (1):
    string System.Web.HttpUtility.UrlDecoder.GetString()
Sanitizers   (8):
    string System.Web.HttpUtility.HtmlEncode(string)
    string System.Web.HttpUtility.UrlEncodeSpaces(string)
    string System.Web.HttpServerUtility.UrlDecode(string)
    string System.Web.HttpUtility.UrlEncode(string, Encoding)
    string System.Web.HttpUtility.UrlEncode(string)
    string System.Web.HttpServerUtility.UrlEncode(string)
    string System.Web.HttpUtility.UrlDecodestringFromstringInternal...
    string System.Web.HttpUtility.UrlDecode(string, Encoding)
Sinks        (4):
    void System.Web.HttpResponse.WriteFile(string)
    void System.Web.HttpRequest.set_QuerystringText(string)
    void System.IO.TextWriter.Write(string)
    void System.Web.HttpResponse.Redirect(string)

```

---

Figure 19: Specification inferred in Example 6.

gone through, as extracted from the code revision repository. It is clear from the table that even arriving at the default specification for CAT.NET, as incomplete as it is, took a pretty significant number of source revisions. Moreover, the initial specification is also fairly large, consisting of a total of 111 methods.

Type	Count	Revisions
Sources	27	11
Sinks	77	10
Sanitizers	7	2

To provide a better metric for the scale of the benchmarks relevant for CAT.NET and MERLIN analyses, Figure 12 provides statistics on the sizes of the propagation graph  $G$  computed by MERLIN, and the factor graph  $F$  constructed in the process of constraint generation. With propagation graphs containing thousands of nodes, it is not surprising that we had to develop a polynomial approximation, as Section 5 describes.

Figure 20: Statistics for the out-of-the box specification that comes with CAT.NET.

## 6.2 Merlin Findings

Figure 13 provides information about the specifications discovered by MERLIN. Columns 2–16 provide information about how many correct and false positive items in each specification category has been found. Note that in addition to “good” and “bad” specifications, as indicated by  $\checkmark$  and  $\times$ , we also have a “maybe” column denoted by  $?$ . This is because often what constitutes a good specification is open to interpretation. Even in consultations with CAT.NET developers we found many cases where the classification of a particular piece of the specification is not clear-cut.

Figure 14 summarizes information about the security vulnerabilities we find based on both the initial and the post-MERLIN specifications. Columns 2–10 show the number of vulnerabilities based on the original specification, newly found vulnerabilities, and the number of newly found vulnerabilities that are in fact false positives. For the post-MERLIN part of the table we also report the number of false positives *eliminated* with the MERLIN specification, 13 in total, as caused by the newly discovered sanitizers. As with many other static analysis tools, false positives is one of the primary complaints about CAT.NET in practice. As can be seen from Figure 14 (the “-” column), MERLIN helps reduce the overall false positive rate from 48% to 33% (the latter computed as  $(43-13)/89$ ).

**Example 5** Function `CreateQueryLink` in Figure 18 is taken from the *Software Catalog* benchmark. The return result of this function is passed

into a known cross-site redirection sink not shown here for brevity.

- The paths that go through `request.Url.AbsolutePath` and `request.QueryString` on lines 6 and 15 are correctly identified vulnerabilities.
- CAT.NET flags the path that passes through function `QueryStringParser.Parse` on line 18 as a vulnerability. However, with MERLIN, `AntiXss.UrlEncode` is correctly determined to be a sanitizer, eliminating this false positives. With MERLIN, we eliminate all 6 false positives in this benchmark.
- CAT.NET addressed explicit information flow only and will not flag the fact that there is a control dependency on line 13 because tainted value `request.QueryString` is used within a conditional.

This short function illustrates many tricky issues with explicit information flow analyses as well as the danger of unrestricted manipulation of tainted data as strings.  $\square$

Note that while we start with the CAT.NET specification characterized in Figure 20, MERLIN can even infer specification entirely *without* an initial specification purely based on the structure of the propagation graph.

**Example 6** Consider a short program fragment written in C# below consisting of two event handlers.

```
protected void TextChanged(object sender, EventArgs e) {
    string str = Request.QueryString["name"];
    string str2 = HttpUtility.HtmlEncode(str);
    Response.Write(str2);
}

protected void ButtonClicked(object sender, EventArgs e) {
    string str = Request.UrlReferrer.AbsolutePath;
    string str2 = HttpUtility.UrlEncode(str);
    Response.Redirect(str2);
}
```

When run with *no initial specification at all*, MERLIN is able to infer a small but absolutely correct specification consisting of 13 elements, as shown in Figure 19. Starting with even a small specification such as the one above, MERLIN is able to successfully infer increasingly larger specifications that fill many gaps in the original CAT.NET specification.  $\square$

### 6.3 Running Times

Finally, Figure 15 provides information about running time of the various MERLIN components, measured in seconds. The experiments were con-

ducted on a 3 GHz Pentium Dual Core Windows XP SP2 machine equipped with 4 GB of memory. Overall, in part due to the approximation described in Section 5, our analysis scales quite well, with none of the benchmarks taking over four minutes to analyze. Given that CAT.NET is generally run once a day or less frequently, these running times are more than acceptable.

## 7 Related Work

Related work falls into the following broad categories discussed in turn: securing Web applications, mining specifications.

### 7.1 Securing Web Applications

There has been much interest in static and runtime protection techniques to improve the security of Web applications. Static analysis allows the developer to avoid issues such as cross-site scripting before the application goes into operation. Runtime analysis allows exploit prevention and recovery during the operation of an application. The WebSSARI project pioneered this line of research [9], by combining static and dynamic analysis for PHP programs. Several projects that came after WebSSARI improve on the quality of static analysis for PHP [10, 33].

The Griffin project proposes a scalable and precise sound static and runtime analysis techniques for finding security vulnerabilities in large Java applications [16, 19]. Several other runtime systems for taint tracking have been proposed as well, including Haldar *et al.* for Java [2, 7] and Pietraszek *et al.* [24] and Nguyen-Tuong *et al.* for PHP [22].

Several tools have been built to detect information flow vulnerabilities in programs [1, 5, 23]. All these tools without exception require a specification of information flow. Our work infers such specifications.

### 7.2 Mining Specifications

We are aware of several efforts to infer security specifications automatically. Tan *et al.*'s [28] AutoISES tool uses static analysis to automatically check if access to sensitive data is suitably guarded by an appropriate authorization check. Similar to our goals, AutoISES automatically infers specifications for this problem. Their core idea is to identify security which data structures are security-sensitive, based on known authorization checks that are present in the code, and check all occurrences of accesses to these data structures to check if they are guarded. Thus, if a security-sensitive data structure is

protected in most accesses, and inadvertently accessed without protection in a few places, they are able to detect such errors. Ganapathy *et al.* used concept analysis to find fingerprints of security-sensitive operations [6]. Unlike these works, we solve a different problem —one of inferring specifications for explicit information flow.

A number of projects have addressed inferring specifications outside the context of security. For a general survey of specification mining techniques, the reader is referred to Perracotta [34], DynaMine [17], and Weimer *et al.* [31]. We mention several efforts below. Engler *et al.* [4] infer specifications from code by seeking rules that involve action pairs: `malloc` paired with `free`, `lock` paired with `unlock`, etc. Li and Zhou [14] and Livshits and Zimmerman [17] look at more general patterns involving action pairs by combining data mining techniques as well as sophisticated pointer analyses. Whaley *et al.* [32] considers inference of interface specifications for Java method calls using static analysis. Jagannathan *et al.* [25] use data mining techniques for inference of method preconditions in complex software systems. The preconditions might incorporate data-flow as well as control-flow properties.

Kremenek *et al.* [11] use probabilistic inference to classify functions that allocate and deallocate resources in programs. While similar in spirit to our work, inference of information flow specifications appears to be a more complex problem than inference of allocation and deallocation routines in C code because there are different kinds of classifiers — sources, sinks, and sanitizers at play. Furthermore, the wrapper avoidance and sanitizers minimization constraints do not have analogs in the allocator-deallocator inference. Unlike Kremenek *et al.* [11] we use the theory of probabilistic refinement to *formally characterize* the triple approximation we have made for the purpose of scaling.

## 8 Conclusions

The growing importance of explicit information flow is evidenced by the abundance of analysis tools for information flow tracking and violation detection at the level of the language, runtime, operating system, and hardware [2, 3, 7, 9–12, 16, 18, 19, 22, 27, 29, 33, 36]. Ultimately, all these approaches need a specification.

In this paper we have presented MERLIN, a novel algorithm that infers explicit information flow specifications from programs. MERLIN derives a system of probabilistic constraints based on interprocedural data flow in the



program, and computes specifications using probabilistic inference.

In order to scale to large programs, we approximate an exponential number of probabilistic constraints by a cubic number of triple constraints, showing that that path-based constraint system is a refinement of the triple-based constraint system. This ensures that, for an given threshold, every solution admitted by the approximated triple system is also admitted by the path system (for the same threshold). Though this connection gives formal grounding to our approximation, it does not say anything about the precision of the results that can be obtained; such an assessment is obtained empirically by evaluating the quality of the specification inferred for large applications, the number of new vulnerabilities discovered, and the number of false positives removed. Based on our observations about large Web applications, we added extra constraints to the triple system (constraints **C2**, **C3**, **C4**, and **C5** in Figure 3) to enhance the quality of the results.

With these extra constraints, our empirical results convincingly demonstrate that our model indeed achieves good precision. In our experiments with 10 large Web applications written in .NET, MERLIN finds a total of 167 new confirmed specifications, which result in a total of 302 newly discovered vulnerabilities across the 10 benchmarks. Equally importantly, MERLIN-inferred specifications also result in 13 false positives being removed. As a result of new findings and eliminating false positives, the *final* false positive rate for CAT.NET *after* MERLIN in our experiments is under 1%.

## References

- [1] CAT.NET. Reference withheld to ensure anonymous submission.
- [2] Deepak Chandra and Michael Franz. Fine-grained information flow analysis and enforcement in a java virtual machine. *Computer Security Applications Conference*, pages 463–475, 10-14 Dec. 2007.
- [3] Michael Dalton, Hari Kannan, and Christos Kozyrakis. Raksha: a flexible information flow architecture for software security. In *ISCA*, pages 482–493, 2007.
- [4] Dawson R. Engler, David Yu Chen, and Andy Chou. Bugs as inconsistent behavior: A general approach to inferring errors in systems code. In *In Proceedings of ACM Symposium on Operating Systems Principles*, pages 57–72, 2001.
- [5] Fortify. Fortify code analyzer. <http://www.ouncelabs.com/>, 2008.

- [6] Vinod Ganapathy, David King, Trent Jaeger, and Somesh Jha. Mining security-sensitive operations in legacy code using concept analysis. In *Proceedings of the International Conference on Software Engineering*, pages 458–467, 2007.
- [7] Vivek Haldar, Deepak Chandra, and Michael Franz. Dynamic taint propagation for Java. In *Proceedings of the Annual Computer Security Applications Conference*, pages 303–311, December 2005.
- [8] C. A. R. Hoare. An axiomatic basis for computer programming. *Communications of the ACM*, 12:576—583, October 1969.
- [9] Yao-Wen Huang, Fang Yu, Christian Hang, Chung-Hung Tsai, Der-Tsai Lee, and Sy-Yen Kuo. Securing Web application code by static analysis and runtime protection. In *Proceedings of the Conference on World Wide Web*, pages 40–52, May 2004.
- [10] Nenad Jovanovic, Christopher Kruegel, and Engin Kirda. Pixy: a static analysis tool for detecting Web application vulnerabilities (short paper). In *Proceedings of the Symposium on Security and Privacy*, May 2006.
- [11] Ted Kremenek, Paul Twohey, Godmar Back, Andrew Y. Ng, and Dawson R. Engler. From uncertainty to belief: Inferring the specification within. In *Symposium on Operating Systems Design and Implementation*, pages 161–176, November 2006.
- [12] Maxwell Krohn, Alexander Yip, Micah Brodsky, Natan Cliffer, M. Frans Kaashoek, Eddie Kohler, and Robert Morris. Information flow control for standard os abstractions. In *Proceedings of Symposium on Operating Systems Principles*, pages 321–334, 2007.
- [13] F. R. Kschischang, B. J. Frey, and H. A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, 47(2):498–519, 2001.
- [14] Z. Li and Y. Zhou. Pr-miner: Automatically extracting implicit programming rules and detecting violations in large software code. In *ICSE*, 2005.
- [15] Benjamin Livshits. *Improving Software Security with Precise Static and Runtime Analysis*. PhD thesis, Stanford University, Stanford, California, 2006.

- [16] Benjamin Livshits and Monica S. Lam. Finding security errors in Java programs with static analysis. In *Proceedings of the Usenix Security Symposium*, pages 271–286, August 2005.
- [17] Benjamin Livshits and Thomas Zimmermann. DynaMine: Finding common error patterns by mining software revision histories. In *Proceedings of the International Symposium on the Foundations of Software Engineering*, pages 296–305, September 2005.
- [18] Michael Martin, Benjamin Livshits, and Monica S. Lam. Finding application errors and security vulnerabilities using PQL: a program query language. In *Proceedings of the Conference on Object-Oriented Programming, Systems, Languages, and Applications*, October 2005.
- [19] Michael Martin, Benjamin Livshits, and Monica S. Lam. SecuriFly: Runtime vulnerability protection for Web applications. Technical report, Stanford University, October 2006.
- [20] A. McIver and C. Morgan. *Abstraction, Refinement and Proof of Probabilistic Systems*. Springer, 2004.
- [21] A. McIver and C. Morgan. Abstraction and refinement in probabilistic systems. *ACM SIGMETRICS Performance Evaluation Review*, 32:41–47, March 2005.
- [22] Anh Nguyen-Tuong, Salvatore Guarnieri, Doug Greene, Jeff Shirley, and David Evans. Automatically hardening Web applications using precise tainting. In *Proceedings of the IFIP International Information Security Conference*, June 2005.
- [23] OunceLabs, Inc. Ounce. <http://www.ouncelabs.com/>, 2008.
- [24] Tadeusz Pietraszek and Chris Vanden Berghe. Defending against injection attacks through context-sensitive string evaluation. In *Proceedings of the Recent Advances in Intrusion Detection*, September 2005.
- [25] M. K. Ramanathan, A. Grama, and S. Jagannathan. Static specification inference using predicate mining. In *PLDI*, 2007.
- [26] Andrei Sabelfeld and Andrew Myers. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications*, 21(1):5–19, January 2003.
- [27] Zhendong Su and Gary Wassermann. The essence of command injection attacks in web applications. In *Proceedings of POPL*, 2006.

- [28] Lin Tan, Xiaolan Zhang, Xiao Ma, Weiwei Xiong, and Yuanyuan Zhou. AutoISES: Automatically inferring security specifications and detecting violations. In *Proceedings of the USENIX Security Symposium*, August 2008.
- [29] Steve Vandebogart, Petros Efstathopoulos, Eddie Kohler, Maxwell Krohn, Cliff Frey, David Ziegler, Frans Kaashoek, Robert Morris, and David Mazières. Labels and event processes in the Asbestos operating system. *ACM Trans. Comput. Syst.*, 25(4):11, 2007.
- [30] Larry Wall. Perl security. <http://search.cpan.org/dist/perl/pod/perlsec.pod>.
- [31] Westley Weimer and George C. Necula. Mining temporal specifications for error detection. In *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 461–476, 2005.
- [32] J. Whaley, M. Martin, and M. Lam. Automatic extraction of object-oriented component interfaces. In *ISSTA*, 2002.
- [33] Yichen Xie and Alex Aiken. Static detection of security vulnerabilities in scripting languages. In *Proceedings of the Usenix Security Symposium*, pages 271–286, August 2006.
- [34] Jinlin Yang and David Evans. Perracotta: mining temporal API rules from imperfect traces. In *Proceedings of the International Conference on Software Engineering*, pages 282–291, 2006.
- [35] Jonathan S. Yedidia, William T. Freeman, and Yair Weiss. Understanding belief propagation and its generalizations. *Exploring Artificial Intelligence in the New Millennium*, pages 239–269, 2003.
- [36] Nikolai Zeldovich, S. Boyd-Wickizer, Eddie Kohler, and David Mazires. Making information flow explicit in HiStar. In *Proceedings of the Symposium on Operating Systems Design and Implementation*, pages 263–278, 2006.