# Technical Report: Maybe Tainted Data, Theory and a Case Study

Christian Skalka [a,*], Sepehr Amir-Mohammadian [b], and Samuel Clark [a]

[a] *University of Vermont, Department of Computer Science*
*E-mails: ceskalka@uvm.edu, samuel.clark@uvm.edu*
[b] *University of the Pacific, Department of Computer Science*
*E-mail: samirmohammadian@pacific.edu*

**Abstract.** Dynamic taint analysis is often used as a defense against low-integrity data in applications with untrusted user interfaces. An important example is defense against XSS and injection attacks in programs with web interfaces. Data sanitization is commonly used in this context, and can be treated as a precondition for endorsement in a dynamic integrity taint analysis. However, sanitization is often incomplete in practice. We develop a model of dynamic integrity taint analysis for Java that addresses imperfect sanitization with an in-depth approach. To avoid false positives, results of sanitization are endorsed for access control (aka prospective security), but are tracked and logged for auditing and accountability (aka retrospective security).

We show how this heterogeneous prospective/retrospective mechanism can be specified as a uniform policy, separate from code. We then use this policy to establish correctness conditions for a program rewriting algorithm that instruments code for the analysis. These conditions synergize our previous work on the semantics of audit logging with *explicit integrity* which is an analogue of noninterference for taint analysis. A technical contribution of our work is the extension of explicit integrity to a high-level functional language setting with structured data, vs. previous systems that only address low level languages with unstructured data. Our approach considers endorsement which is crucial to address sanitization. An implementation of our rewriting algorithm is presented that hardens the OpenMRS medical records software system with in-depth taint analysis, along with an empirical evaluation of the overhead imposed by instrumentation. Our results show that this instrumentation is practical.

Keywords: Auditing, Dynamic taint analysis, Program rewriting

## 1. Introduction

Dynamic taint analysis implements a "direct" or "explicit" information flow analysis to support a variety of security mechanisms [1]. Similar to information flow, taint analysis can be used to support either confidentiality or integrity properties. An important application of integrity taint analysis is to prevent the execution of security sensitive operations on untrusted data, in particular to combat cross-site scripting (XSS) and SQL injection attacks in web applications [2]. Any untrusted user input is marked as tainted, and then taint is tracked and propagated through data flow to ensure that tainted data is not used by security sensitive operations.

Of course, since web applications aim to be interactive, user input is needed for certain security sensitive operations such as database calls. To combat this, *sanitization* is commonly applied in practice to analyze and possibly modify data. From a taint analysis perspective, sanitization is a precondition for integrity endorsement, i.e. subsequently viewing sanitization results as high integrity data. However, while

sanitization is usually endorsed as "perfect" by taint analysis, in fact it is not. Indeed, previous work has identified a number of flaws in existing sanitizers in a variety of applications [2, 3]. The work here is in fact inspired by discovery of an XSS attack vector in the OpenMRS medical records software systems due to incomplete sanitization, discussed below in Section 1.1. We call such incomplete sanitizers *partially trusted* or *imperfect* throughout the paper.

Thus, a main challenge we address is how to mitigate imperfect sanitization in taint analysis. An important feature of our approach is an *in-depth* [4] security policy, that combines the typical blocking (prospective) behavior of taint-based access control with audit logging (retrospective) features. In the presence of imperfect sanitization, this allows false positives to be avoided, while still providing retrospective security measures via audit logs in case of attacks that leverage this imperfection. We are concerned with both efficiency and correctness– we develop a language model intended to capture the essence of Phosphor [5, 6], an existing Java taint analysis system with empirically demonstrated efficiency. We define a source code language based on Featherweight Java (FJ) [7], and propose a rewriting algorithm that instruments code with support for in-depth integrity taint analysis in the presence of partially trusted sanitization. This model forms the basis for development of a real implementation of in-depth taint analysis for OpenMRS, also discussed in this paper.

A contribution of our approach is a uniform expression of an in-depth security policy that combines prospective (taint analysis) and retrospective (audit logging) policy features, and proof that our rewriting algorithm enforces this policy. To characterize retrospective correctness we leverage our previous work on the the semantics of retrospective security [8]. To characterize prospective correctness, we aim to go deeper than operational definitions [1, 9] and characterize correctness as a true security property. To this end we propose a semantic framework called *explicit integrity*, which is an extension of explicit secrecy [10] to a high-level (Java) language model with structured data. Both explicit secrecy and integrity are defined independently of language-level instrumentation, and (like e.g. information flow noninterference) are *hyperproperties* [11]. Intuitively, if a program satisfies explicit integrity, then it is guaranteed that data from low-integrity sources does not *directly* (i.e. *explicitly*) flow into high-integrity sinks. Furthermore, we consider the variant explicit integrity *modulo endorsement*, since endorsement is necessary in the taint analysis to accurately reflect the results of sanitization.

### 1.1. Vulnerability and Countermeasures

While our work is based on formal foundations it is inspired by practical concerns, in particular a security flaw in OpenMRS. This flaw allows an attacker to launch persistent XSS attacks[1]. When a web-based software receives and stores user input without proper sanitization, and later retrieves this information for (other) users, persistent XSS attacks could take place.

OpenMRS uses a set of validators to enforce expected data formats by implementation of the `Validator` interface (e.g., `PersonNameValidator`, `VisitTypeValidator`, etc.). For some of these classes the implementation is strict enough to reject script tags by enforcing data to match a particular regular expression, e.g., `PersonNameValidator`. However, `VisitTypeValidator` lacks such restriction and only checks for object fields to avoid being null, empty or whitespace, and their lengths to be correct. Thus the corresponding webpage that receives user inputs to construct `VisitType` objects (named `VisitTypeForm.jsp`) is generally not able to perform proper sanitization through the invocation of the validator implemented by `VisitTypeValidator`. A `VisitType`

---

[1]We have responsibly disclosed the vulnerabilities we have found in OpenMRS (version 2.4, released 7/15/2016, and the preceding versions) to the OpenMRS developing community. We discuss one particular case here.

object is then stored in the MySQL database, and could be retrieved later based on user request. For instance, `VisitTypeList.jsp` queries the database for all defined `VisitType` objects, and sends `VisitType` names and descriptions to the client side. Therefore, the attacker can easily inject scripts as part of `VisitType` name and/or description, and the constructed object would be stored in the database and possibly in a later stage retrieved and executed in the victim's client environment.

Integrity taint tracking is a well-recognized solution against these sorts of attacks. In our example, using taint analysis the tainted `VisitType` object would be prevented from retrieval and execution. The addition of sanitization methods would also be an obvious step, and commensurate with an integrity taint analysis approach– sanitized objects would be endorsed for the purposes of prospective security. However, many attack scenarios demonstrate degradation of taint tracking effectiveness due to unsound or incomplete input sanitization [2, 3].

To support integrity taint analysis in the presence of incomplete sanitization for legacy code, we propose a program rewriting approach, which is applicable to systems such as OpenMRS. Our program rewriting algorithm takes as input a heterogeneous (prospective and retrospective) taint analysis policy specification and input code, and instruments the code to support the policy. The policy allows user specification of taint sources, secure sinks, and sanitizers. A distinct feature of our system is that results of sanitization are considered "maybe tainted" data, which are allowed to flow into security sensitive operations but in such cases are entered in a log to support auditing and accountability.

### 1.2. The Security and Threat Model

The security problem we consider is about the integrity of data being passed to security sensitive operations (*SSOs*). An important example is a string entered by an untrusted user that is passed to a database method for parsing and execution as a SQL command. The security mechanism should guarantee that low-integrity data cannot be passed to *SSOs* without previous sanitization.

In contrast to standard information flow which is concerned with both direct (aka explicit) and indirect (aka implicit) flows, taint analysis is only concerned with direct flow. Direct flows transfer data directly between variables, e.g., $n_1$ and $n_2$ directly affect the result of $n_1 + n_2$. Indirect flows are realized when data can affect the result of code dispatch– the standard example is a conditional expression **if** $v$ **then** $e_1$ **else** $e_2$ where the data $v$ indirectly affects the valuation of the expression by guarding dispatch.

More precisely, we posit that top-level programs $\mathfrak{p}$ in this security setting are parameterized by a low integrity data source **a**, and an arbitrary number of secure sinks (*SSOs*) and sanitizers which are specified externally to the program by a security administrator. For simplicity we assume that *SSOs* are unary operations over primitive objects, so there is no question about which argument may be tainted. Since we define a Java based model, each sso or sanitizer is identified as a specific method m in a class C. That is, there exists a set of *Sanitizers* containing class, method pairs C.m which are assumed to return high-integrity data, though they may be passed low-integrity data. Likewise, there exists a set of *SSOs* of the same form, and for brevity we will write $sso(e)$ for a method invocation v.m(e) on some object v where C.m $\in$ *SSO*. As a sanity condition we require *SSOs* $\cap$ *Sanitizers* $= \varnothing$. For simplicity of our formal presentation we assume that only one tainted source will exist. Explicit integrity, as a high-level property, is instantiated for this model.

We assume that our program rewriting algorithm is trusted. Input code is trusted to be not malicious, though it may contain errors. We note that this assumption is important for application of taint analysis that disregards indirect flows, since there is confidence that the latter will not be exploited (even accidentally) as a side-channel attack vector by non-malicious code. We assume that untrusted data sources

provide low integrity data, though in this work we only consider tainted "static" values, e.g., strings, not tainted code that may be run as part of the main program execution. However, the latter does not preclude hardening against XSS or injection attacks in practice, if we consider an evaluation method to be an sso.

### 1.3. Overview and Main Technical Results of this Paper

The technical development of the paper proceeds as follows. In Section 2 we describe a formal semantics of auditing, and the conditions for correctness of audit rewriting algorithms. That is, we define what it means for a program instrumentation to correctly log information. In Section 2.1, we introduce information algebra [12] as the basis of our model for correct audit log generation. We characterize logging specifications and correctness conditions for audit logs in a high-level manner using information algebra, and show how information elements and operations can be instantiated using first order logic.

In Section 3 we develop a source language model based on featherweight Java (Section 3.1), called FJ. We show how to logically specify an in-depth taint analysis policy separately from code in Section 4 via safety property and logging specifications. In Section 3.2 we develop a target language model $FJ_{taint}$ with instrumentation for operationally enforcing an in-depth taint analysis, which we show to be correct according to our formal condition in Section 4.2, with our main result being Theorem 4.1.

While Theorem 4.1 establishes correctness conditions for information in audit logs in an operational sense, Section 5 focuses on the high level security property of dynamic integrity taint analysis. In Section 5.4, we show that our enforcement mechanism satisfies the hyperproperty of explicit integrity modulo endorsement (Theorem 5.1).

In Section 6 we discuss our implementation of the in-depth taint analysis specification presented in Section 4 for the OpenMRS medical records system. In Section 7 we conclude the paper.

### 1.4. Related Work

Taint analysis is an established solution to enforce confidentiality and integrity policies through direct data flow control. Various systems have been proposed for both low and high level level languages. Our policy language and semantics are based on a well-developed formal foundation, where we interpret Horn clause logic as an instance of information algebra [13] in order to specify and interpret retrospective policies. The work presented in this paper supersedes a previous presentation [14]– in the current paper we extend our language model, provide more rigorous proofs of correctness of policy enforcement, consider the hyperproperty of taint analysis in a model of Java, and report on a prototype implementation.

Schwartz et al. [1] define a general model for runtime enforcement of policies using taint tracking for an intermediate language. In Livshits [9], taint analysis is expressed as part of operational semantics, similar to Schwartz et al., and a taxonomy of taint tracking is defined. Livshits et al. [2] propose a solution for a range of vulnerabilities regarding Java-based web applications, including SQL injections, XSS attacks and parameter tampering, and formalize taint propagation including sanitization. The work uses PQL [15] to specify vulnerabilities. However, these works are focused on operational definitions of taint analysis for imperative languages. In contrast we have developed a logical specification of taint analysis for a functional OO language model that is separate from code, and is used to establish correctness of an implementation. Our work also comprises a unique retrospective component to protect against incomplete input sanitization. According to earlier studies [2, 3], incomplete input sanitization makes a variety of applications susceptible to injection attacks.

Another related line of work is focused on the optimization of integrity taint tracking deployment in web-based applications. Sekar [16] proposes a taint tracking mechanism to mitigate injection attacks in web applications. The work focuses on input/output behavior of the application, and proposes a lower-overhead, language-independent and non-intrusive technique that can be deployed to track taint information for web applications by blackbox taint analysis with syntax-aware policies. In our work, however, we propose a deep instrumentation technique to enforce taint propagation in a layered in-depth fashion. Wei et al. [17] attempt to lower the memory overhead of TaintDroid taint tracker [18] for Android applications. The granularity of taint tracking places a significant role in the memory overhead. To this end, TaintDroid trades taint precision for better overhead, e.g., by having a single taint label for an array of elements. Our work reflects a more straightforward object-level taint approach in keeping with existing Java approaches.

Saxena et al. [19] employ static techniques to optimize dynamic taint tracking done by binary instrumentation, through the analysis of registers and stack frames. They observe that it is common for multiple local memory locations and registers to have the same taint value. A single taint tag is used for all such locations. A shadow stack is employed to retain the taint of objects in the stack. Cheng et al. [20] also study the solutions for taint tracking overhead for binary instrumentation. They propose a byte to byte mapping between the main and shadow memory that keeps taint information. Bosman et al. [21] propose a new emulator architecture for the x86 architecture from scratch with the sole purpose of minimizing the instructions needed to propagate taint. Similar to Cheng et al. [20], they use shadow memory to keep taint information, with a fixed offset from user memory space. Zhu et al. [22] track taint for confidentiality and privacy purposes. In case a sensitive input is leaked, the event is either logged, prohibited or replaced by some random value. We have modeled a similar technique for an OO language, through high level logical specification of shadow objects, so that each step of computation is simulated for the corresponding shadow expressions.

Particularly for Java, Chin et al. [23] propose taint tracking of Java web applications in order to prohibit injection attacks. To this end, they focus on strings as user inputs, and analyze the taint in character level. For each string, a separate taint tag is associated with each character of the string, indicating whether that character was derived from untrusted input. The instrumentation is only done on the string-related library classes to record taint information, and methods are modified in order to propagate taint information. Haldar et al. [24] propose an object-level tainting mechanism for Java strings. They study the same classes as the ones in Chin et al. [23], and instrument all methods in these classes that have some string parameters and return a string. Then, the returned value of instrumented method is tainted if at least one of the argument strings is tainted. However, in contrast to our work, *only* strings are endowed with integrity information, whereas all values are assigned integrity labels in our approach. Recently Bodei et al. [25] have proposed a static enforcement mechanism for taint analysis in IoT devices which predicts the propagation of taint in the system according to the flow of control. These previous works lack retrospective features.

Phosphor [5, 6] is an attempt to apply taint tracking more generally in Java, to any primitive type and object class. Phosphor instruments the application and libraries at bytecode level based on a given list of taint source and sink methods. Input sanitizers with endorsement are not directly supported, however. As Phosphor avoids any modifications to the JVM, the instrumented code is still portable. Our work is an attempt to formalize Phosphor in FJ extended with input sanitization and in-depth enforcement. Our larger goal is to develop an implementation of in-depth dynamic integrity analysis for Java by leveraging the existing Phosphor system.

Secure information flow [26] and its interpretation as the well-known hyperproperty [11] of noninterference [27] is challenging to implement in practical settings [28] due to implicit flows. Taint analysis is thus an established solution to enforce confidentiality and integrity policies since it tracks only direct data flow control. Various systems have been proposed for both low and high level level languages. The majority of previous work, however, has been focused on taint analysis policy specification and enforcement (e.g., [1, 9, 17, 29]), rather than capturing the essence of direct information flow which could provide an underlying framework to study numerous taint analysis tools.

Knowledge-based semantics has been introduced by Askarov et al. [30] as a general model for information flow of confidential data, concentrated on cryptographic computations and key release (declassification [31]) and later employed in other data secrecy analyses [32–34]. Schoepe et al. [10] have proposed the semantic notion of correctness for taint tracking that enforces confidentiality policies of direct information flow, called explicit secrecy. To this end, they propose a knowledge-based semantics, influenced by Volpano's weak secrecy [35] and gradual release [30]. Explicit secrecy is defined as a property of a program, where the program execution does not change the explicit knowledge of public user. The authors show that noninterference is not comparable to explicit secrecy. However, rather than restricting the discussion to direct information flow in a low level language, we model a high level OO language with a functional flavor to represent generality of our framework.

Schoepe et al. [36] have recently employed explicit secrecy to study correctness results for dynamic confidentiality taint analysis in a core imperative setting with pointers and I/O, and deployed a Java-based tool, called DroidFace. A recent framework by Balliu et al. [37] attempts to bring together the general information flow and direct flow analyses using a security condition that models indirect flows which are observable by a low confidentiality user.

A counterpart for attacker knowledge in the realm of general flow of information integrity, called attacker power [34], is introduced as the set of low integrity inputs that generate the same observables. In this regard, Askarov et al. [34] use holes in the syntax of program code for injection points, influenced by [38]. However, their attack model is different as the low integrity and low confidentiality user is able to inject program code in the main program, by which she could gain more knowledge. We have tailored attacker power for explicit flows using state transformers, in order to interpret integrity taint analysis.

Birgisson et al. [39] give a unified framework to capture different flavors of integrity, in particular integrity via information flow and via different types of invariance. Similar to other works in this line, they give a simple imperative language with labeled operational semantics in order to enforce integrity policies through communication with a monitor. In contrast, we use program rewriting techniques to enforce policies regarding flow of data integrity, which are applicable to legacy systems.

In addition to formal properties of direct information flow, our formulation of correctness conditions also considers a formalization of audit logging based on our previous work [8], which considered a safety property unrelated to taint analysis. Other authors have recently considered formal characterizations of auditing based on logics of justification [40, 41]. In contrast, we consider a specific security application of auditing in combination with taint analysis where audit logs are "extralinguistic" vestiges of program computation, whereas these related works consider programs that are able to reflect on their own audit trails, which is a distinct theoretical problem.

## 2. Foundations for In-Depth Policy Specification

In this Section we establish formal foundations for a semantics of prospective and retrospective policy features. More specifically, we develop a framework for characterizing the correctness of audit instru-

mentation, though safety properties in the standard style [42] can also be expressed in our framework–thus it is appropriate for uniformly characterizing operational correctness conditions for in-depth integrity taint analysis. This framework was initially developed in previous work [8] where we studied so-called "break-the-glass" policies for medical records software. In that work we justified the generality of our framework and discuss its details at length, here we reiterate the main technical points of the framework to allow a standalone formal presentation.

We leverage ideas from the theory of *information algebra* [12, 43], which is an abstract mathematical framework for information systems. In short, we interpret program traces as information, and logging specifications as functions from traces to information. This separates logging specifications from their implementation in code, and defines exactly the information that should be in an audit log. This in turn establishes correctness conditions for audit logging implementations.

## 2.1. Introduction to Information Algebra

*Information algebra* is an algebraic theory of information where information is seen as a collection of *information elements* with fundamental aggregation and refinement operations. The algebra consists of two domains, an information domain and a query domain. The information domain $\Phi$ is the set of information elements that can be aggregated in order to build more inclusive information elements. The query domain $E$ is a lattice of querying sublanguages in which the partial order relation among these sublanguages represents the granularity of the queries. The information and query domains are left abstract in the general theory– instantiation examples include relational algebra and first order logic as we discuss below. By definition any instantiation must include basic operations for combining information and for focusing on components of information.

**Definition 2.1.** *Any information algebra* $(\Phi, E)$ *includes two basic operators:*

- *Combination* $\otimes : \Phi \times \Phi \to \Phi$*: The operation* $X \otimes Y$ *combines (or, aggregates)* the information in *elements* $X, Y \in \Phi$.
- *Focusing* $\Rightarrow : \Phi \times E \to \Phi$*: The operation* $X^{\Rightarrow S}$ *isolates the elements of* $X \in \Phi$ *that are relevant to a sublanguage* $S \in E$*, i.e. the subpart of X specified by S.*

Using the combination operator we can define a partial order relation on $\Phi$ to compare the information contained in the elements of $\Phi$. A partial ordering is induced on $\Phi$ by the so-called *information ordering* relation $\leqslant$, where intuitively for $X, Y \in \Phi$ we have $X \leqslant Y$ iff $Y$ contains at least as much information as $X$, though its precise meaning depends on the particular algebra.

**Definition 2.2.** *X is contained in Y, denoted as* $X \leqslant Y$*, for all* $X, Y \in \Phi$ *iff* $X \otimes Y = Y$.

**Definition 2.3.** *We say that X and Y are* information equivalent, *and write* $X = Y$*, iff* $X \leqslant Y$ *and* $Y \leqslant X$.

For a more detailed account of information algebra, the reader is referred to a definitive survey paper [43].

### 2.1.1. Illustrative Example: Relational Algebras

Relational algebra is a well-recognized instance of information algebra. Let $\mathcal{A}$ denote the set of *attributes*, $\mathcal{A}_i \subseteq \mathcal{A}$ for $i \in \{1, 2, 3\}$, $\mathcal{A}_2 \subseteq \mathcal{A}_1$, and assume that $\mathcal{A}_1 = \{a_1, ..., a_n\}$. Each tuple $((a_1 : x_1), \cdots, (a_n : x_n))$ can be formulated as a function $f : \mathcal{A}_1 \to \{x_1, ..., x_n\}$, where $f(a_i) = x_i$.

Function $f[\mathcal{A}_2] : \mathcal{A}_2 \rightarrow \{x_1, ..., x_n\}$ is the *restriction* of $f$ to $\mathcal{A}_2$, defined as $f[\mathcal{A}_2](a) = f(a)$, for all $a \in \mathcal{A}_2$. A *relation* $R$ over $\mathcal{A}_1$ is a set of functions $f$ defined on a specific set of attributes $\mathcal{A}_1$. Then, the *projection* of $R$ on $\mathcal{A}_2$ is defined as $\pi_{\mathcal{A}_2}(R) = \{f[\mathcal{A}_2] \mid f \in R\}$. The *natural join* of relation $R$ over $\mathcal{A}_1$ and $R'$ over $\mathcal{A}_3$ is defined as $R \bowtie R' = \{f \mid dom(f) = \mathcal{A}_1 \cup \mathcal{A}_3, f[\mathcal{A}_1] \in R, f[\mathcal{A}_3] \in R'\}$.

*Instantiation.*    Let $\mathbb{R}$ be the universe of all relations $R$. Then, $(\mathbb{R}, \mathcal{P}(\mathcal{A}))$ is an information algebra with following definitions for combination and focusing:

$$R \otimes R' \triangleq R \bowtie R' \qquad\qquad R^{\Rightarrow \mathcal{A}_1} \triangleq \pi_{\mathcal{A}_1}(R)$$

## 2.2. A General Model for Logging Specifications

Following [42], an *execution trace* $\tau = \kappa_0 \kappa_1 \kappa_2 \ldots$ is a possibly infinite sequence of configurations $\kappa$ that describe the state of an executing program. We deliberately leave configurations abstract, but examples abound and we explore a specific instantiation for FJ-based calculus in Section 3. Note that an execution trace $\tau$ may represent the partial execution of a program, i.e. the trace $\tau$ may be extended with additional configurations as the program continues execution. We use metavariables $\tau$ and $\sigma$ to range over traces, and use $\varnothing$ to denote an empty trace.

We assume given a function $\lfloor \cdot \rfloor$ that is an injective mapping from traces to $\Phi$. This mapping *interprets a given trace as information*, where the injective requirement ensures that information is not lost in the interpretation. For example, if $\sigma$ is a proper prefix of $\tau$ and thus contains strictly less information, then formally $\lfloor \sigma \rfloor \leqslant \lfloor \tau \rfloor$. We intentionally leave both $\Phi$ and $\lfloor \cdot \rfloor$ underspecified for generality, though application of our formalism to a particular logging implementation requires instantiation of them.

We let *LS* range over *logging specifications*, which are functions from traces to $\Phi$. As for $\Phi$ and $\lfloor \cdot \rfloor$, we intentionally leave the language of specifications abstract, but consider a particular instantiation in Section 2.6. Intuitively, $LS(\tau)$ denotes the information that should be recorded in an audit log during the execution of $\tau$ given specification *LS*, regardless of whether $\tau$ actually records any log information, correctly or incorrectly. We call this the semantics of the logging specification *LS*.

We assume that auditing is implementable, requiring at least that all conditions for logging any piece of information must be met in a finite amount of time. As we will show, this restriction implies that correct logging instrumentation is a safety property [42].

**Definition 2.4.** *We require of any logging specification LS that for all traces $\tau$ and information $X \leqslant LS(\tau)$, there exists a finite prefix $\sigma$ of $\tau$ such that $X \leqslant LS(\sigma)$.*

It is crucial to observe that some logging specifications may *add* information not contained in traces to the auditing process. Security information not relevant to program execution (such as ACLs), interpretation of event data (statistical or otherwise), etc., may be added by the logging specification. For example, in the OpenMRS system [44], logging of sensitive operations includes a human-understandable "type" designation which is not used by any other code. Thus, given a trace $\tau$ and logging specification *LS*, it is *not* necessarily the case that $LS(\tau) \leqslant \lfloor \tau \rfloor$. Audit logging is not just a filtering of program events.

## 2.3. Correctness Conditions for Audit Logs

A logging specification defines what information should be contained in an audit log. In this section we develop formal notions of *soundness* and *completeness* as audit log correctness conditions. We use

metavariable $\mathbb{L}$ to range over audit logs. Again, we intentionally leave the language of audit logs unspecified, but assume that the function $\lfloor \cdot \rfloor$ is extended to audit logs, i.e. $\lfloor \cdot \rfloor$ is an injective mapping from audit logs to $\Phi$. Intuitively, $\lfloor \mathbb{L} \rfloor$ denotes the information in $\mathbb{L}$, interpreted as an element of $\Phi$.

An audit log $\mathbb{L}$ is sound with respect to a logging specification *LS* and trace $\tau$ if the log information is contained in $LS(\tau)$. Similarly, an audit log is complete with respect to a logging specification if it contains all of the information in the logging specification's semantics. Crucially, both definitions are independent of the implementation details that generate $\mathbb{L}$.

**Definition 2.5.** *Audit log* $\mathbb{L}$ *is* sound *with respect to logging specification LS and execution trace $\tau$ iff* $\lfloor \mathbb{L} \rfloor \leqslant LS(\tau)$.

**Definition 2.6.** *Audit log* $\mathbb{L}$ *is* complete *with respect to logging specification LS and execution trace $\tau$ iff* $LS(\tau) \leqslant \lfloor \mathbb{L} \rfloor$.

*2.4. Correct Logging Instrumentation is a Safety Property*

In case program executions generate audit logs, we write $\tau \rightsquigarrow \mathbb{L}$ to mean that trace $\tau$ generates $\mathbb{L}$, i.e. $\tau = \kappa_0 \ldots \kappa_n$ and $logof(\kappa_n) = \mathbb{L}$ where $logof(\kappa)$ denotes the audit log in configuration $\kappa$, i.e. the residual log after execution of the full trace. Ideally, information that *should* be added to an audit log, *is* added to an audit log, immediately as it becomes available. This ideal is formalized as follows.

**Definition 2.7.** *For all logging specifications LS, the trace $\tau$ is* ideally instrumented *for LS iff for all finite prefixes $\sigma$ of $\tau$ we have $\sigma \rightsquigarrow \mathbb{L}$ where $\mathbb{L}$ is sound and complete with respect to LS and $\sigma$.*

We observe that the restriction imposed on logging specifications by Definition 2.4, implies that ideal instrumentation of any logging specification is a safety property in the sense defined by Schneider [42].

**Theorem 2.1.** *For all logging specifications LS, the set of ideally instrumented traces is a safety property.*

**Proof.** If $\tau$ is ideally instrumented for *LS*, then it is prefix-closed by definition. Furthermore, if $\tau$ is not ideally instrumented for *LS*, then it will definitely be rejected in a finite amount of time, since any information in $LS(\tau)$ is encountered after execution of a finite prefix $\sigma$ of $\tau$ by Definition 2.4. These two facts obtain the result. □

This result implies that e.g. edit automata can be used to enforce instrumentation of logging specifications [45]. However, theory related to safety properties and their enforcement by execution monitors [42, 46] do not provide an adequate semantic foundation for audit log generation, nor an account of soundness and completeness of audit logs.

*2.5. Implementing Logging Specifications with Program Rewriting*

The above-defined correctness conditions for audit logs provide a foundation on which to establish correctness of logging implementations. Here we consider program rewriting approaches. Since rewriting concerns specific languages, we introduce an abstract notion of programs **p** with an operational semantics that can produce a trace. We write $\mathbf{p} \Downarrow \sigma$ iff program **p** can produce execution trace $\tau$, either deterministically or non-deterministically, and $\sigma$ is a *finite* prefix of $\tau$.

A rewriting algorithm $\mathcal{R}$ is a (partial) function that takes a program $\mathbf{p}$ in a source language and a logging specification *LS* and produces a new program, $\mathcal{R}(\mathbf{p}, LS)$, in a target language.[2] The intent is that the target program is the result of instrumenting $\mathbf{p}$ to produce an audit log appropriate for the logging specification *LS*. A rewriting algorithm may be partial, in particular because it may only be intended to work for a specific set of logging specifications.

Ideally, a rewriting algorithm should preserve the semantics of the program it instruments. That is, $\mathcal{R}$ is semantics-preserving if the rewritten program simulates the semantics of the source code, modulo logging steps. We assume given a correspondence relation $\cong$ on execution traces. A coherent definition of correspondence should be similar to a bisimulation, but it is not necessarily symmetric nor a bisimulation, since the instrumented target program may be in a different language than the source program. We deliberately leave the correspondence relation underspecified, as its definition will depend on the instantiation of the model. Possible definitions are that traces produce the same final value, or that traces when restricted to a set of memory locations are equivalent up to stuttering. Furthermore, because rewriting will often add blocking checks for unsafe behaviors (as in the case we will study), semantics preservation is defined up to simulation of sets of program traces that will typically be defined as a safety property. We provide an explicit definition of correspondence for FJ-calculus source and target languages in Section 4.2.

**Definition 2.8.** *Let T be a set of program traces. Rewriting algorithm $\mathcal{R}$ is* semantics preserving up to *T iff for all programs $\mathbf{p}$ and logging specifications LS such that $\mathcal{R}(\mathbf{p}, LS)$ is defined, all of the following hold:*

(1) *For all traces $\tau \in T$ such that $\mathbf{p} \Downarrow \tau$ there exists $\tau'$ with $\tau \cong \tau'$ and $\mathcal{R}(\mathbf{p}, LS) \Downarrow \tau'$.*
(2) *For all traces $\tau$ such that $\mathcal{R}(\mathbf{p}, LS) \Downarrow \tau$ there exists a trace $\tau' \in T$ such that $\tau' \cong \tau$ and $\mathbf{p} \Downarrow \tau'$.*

In addition to preserving program semantics, a correctly rewritten program constructs a log in accordance with the given logging specification. More precisely, if *LS* is a given logging specification and a trace $\tau$ describes execution of a source program, rewriting should produce a program with a trace $\tau'$ that corresponds to $\tau$ (i.e., $\tau \cong \tau'$), where the log $\mathbb{L}$ generated by $\tau'$ contains the same information as $LS(\tau)$, or at least a sound approximation. Some definitions of $\cong$ may allow several target-language traces to correspond to source-language traces (as for example in Section 4.2, our Definition 4.7). In any case, we expect that at least one simulation exists. Hence we write *simlogs*$(\mathbf{p}, \tau)$ to denote a nonempty set of logs $\mathbb{L}$ such that, given source language trace $\tau$ and target program $\mathbf{p}$, there exists some trace $\tau'$ where $\mathbf{p} \Downarrow \tau'$ and $\tau \cong \tau'$ and $\tau' \rightsquigarrow \mathbb{L}$. The name *simlogs* evokes the relation to logs resulting from simulating executions in the target language.

The following definitions then establish correctness conditions for rewriting algorithms. Note that satisfaction of either of these conditions only implies condition (i) of Definition 2.8, not condition (ii), so semantics preservation is an independent condition. Like semantics preservation, we define soundness and completeness with respect to a given set of traces.

**Definition 2.9.** *Let T be a set of traces. Rewriting algorithm $\mathcal{R}$ is* sound up to *T iff for all programs $\mathbf{p}$, logging specifications LS, and finite traces $\tau \in T$ where $\mathbf{p} \Downarrow \tau$, for all $\mathbb{L} \in simlogs(\mathcal{R}(\mathbf{p}, LS), \tau)$ it is the case that $\mathbb{L}$ is sound with respect to LS and $\tau$.*

---

[2]We use metavariable $\mathbf{p}$ to range over programs in either the source or target language; it will be clear from context which language is used.

**Definition 2.10.** *Let $T$ be a set of traces. Rewriting algorithm $\mathcal{R}$ is* complete up to $T$ *iff for all programs* **p**, *logging specifications LS, and finite traces $\tau \in T$ where* $\mathbf{p} \Downarrow \tau$, *for all $\mathbb{L} \in simlogs(\mathcal{R}(\mathbf{p}, LS), \tau)$ it is the case that $\mathbb{L}$ is complete with respect to LS and $\tau$.*

*2.6. A First Order Logic (FOL) Specification Language*

Logics have been used in several well-developed auditing systems [47, 48], for the encoding of both audit logs and queries. FOL in particular is attractive due to readily available implementation support, e.g. Datalog and Prolog. We have shown in previous work that FOL is an information algebra, and useful for e.g. break the glass policy specification [8]. Here we summarize important definitions for the remainder of this paper.

Let Greek letters $\phi$ and $\psi$ range over FOL formulas and let capital letters $X, Y, Z$ range over sets of formulas. We posit a sound and complete proof theory supporting judgements of the form $X \vdash \phi$. In this text we assume without loss of generality a natural deduction proof theory.

Elements of our algebra are sets of formulas closed under logical entailment. Intuitively, given a set of formulas $X$, the closure of $X$ is the set of formulas that are logically entailed by $X$, and thus represents all the information contained in $X$. In spirit, we follow the treatment of sentential logic as an information algebra explored in related foundational work [12], however our definition of closure is syntactic, not semantic.

**Definition 2.11.** *We define a closure operation $C$, and a set $\Phi_{FOL}$ of closed sets of formulas:*

$$C(X) = \{\phi \mid X \vdash \phi\} \qquad \Phi_{FOL} = \{X \mid C(X) = X\}$$

*Note in particular that $C(\varnothing)$ is the set of logical tautologies.*

Let *Preds* be the set of all predicate symbols, and let $S \subseteq Preds$ be a set of predicate symbols. We define *sublanguage $L_S$* to be the set of well-formed formulas over predicate symbols in $S$ (including boolean atoms *true* and *false*, and closed under the usual first-order connectives and binders). We will use sublanguages to define refinement operations in our information algebra. Subset containment induces a lattice structure, denoted $\mathcal{S}$, on the set of all sublanguages, with $\mathcal{F} = L_{Preds}$ as the top element.

Now we can define the focusing and combination operators, which are the fundamental operators of an information algebra. Focusing isolates the component of a closed set of formulas that is in a given sublanguage. Combination closes the union of closed sets of formulas. Intuitively, the focus of a closed set of formulas $X$ to sublanguage $L$ is the refinement of the information in $X$ to the formulas in $L$. The combination of closed sets of formulas $X$ and $Y$ combines the information of each set.

**Definition 2.12.** *Define:*

(1) *Focusing: $X^{\Rightarrow S} = C(X \cap L_S)$ where $X \in \Phi_{FOL}$, $S \subseteq Preds$*
(2) *Combination: $X \otimes Y = C(X \cup Y)$ where $X, Y \in \Phi_{FOL}$*

Properties of the algebra ensure that $\leqslant$ is a partial ordering by defining $X \leqslant Y$ iff $X \otimes Y = Y$, which in the case of our logical formulation means that for all $X, Y \in \Phi_{FOL}$ we have $X \leqslant Y$ iff $X \subseteq Y$, i.e. $\leqslant$ is subset inclusion over closed sets of formulas.

The following Theorem establishes that the construction is an information algebra– for a complete proof the reader is directed to [45].

**Theorem 2.2.** *Structure* $(\Phi_{FOL}, \mathcal{S})$ *with focus operation* $X^{\Rightarrow S}$ *and combination operation* $X \otimes Y$ *forms a domain-free information algebra.*

In addition, to interpret traces and logs as elements of this algebra, i.e. to define the function $\lfloor \cdot \rfloor$, we assume existence of a function $toFOL(\cdot)$ that injectively maps traces and logs to sets of FOL formulas, and then take $\lfloor \cdot \rfloor = C(toFOL(\cdot))$. To define the range of $toFOL(\cdot)$, that is, to specify how trace information will be represented in FOL, we assume the existence of *configuration description predicates P* which are each at least unary. Each configuration description predicate fully describes some element of a configuration $\kappa$, and the first argument is always a natural number $n$, indicating the time at which the configuration occurred. A set of configuration description predicates with the same timestamp describes a configuration, and traces are described by the union of sets describing each configuration in the trace. We will fully define $toFOL(\cdot)$ when we discuss particular source and target languages for program rewriting.

Formally, we define logging specifications in a logic programming style by using combination and focusing. Any logging specification is parameterized by a sublanguage $S$ that identifies the predicate(s) to be resolved and Horn clauses $X$ that define it/them, hence we define a functional *spec* from pairs $(X, S)$ to specifications $LS$, where we use $\lambda$ as a binder for function definitions in the usual manner:

**Definition 2.13.** *The function spec is given a pair* $(X, S)$ *and returns a* FOL *logging specification, i.e. a function from traces to elements of* $\Phi_{FOL}$:

$$spec(X, S) = \lambda\tau.(\lfloor \tau \rfloor \otimes C(X))^{\Rightarrow S}.$$

## 3. Direct Information Flow: Dynamic Integrity Taint Analysis

In this section we present a basic object-oriented calculus as the foundation of our language model. We also show how the in-depth integrity taint analysis model described in Section 1.2 can be specified as a logical property of program traces in this model, independent of program instrumentation. This allows us to define retrospective taint analysis as a logging specification in the style introduced in Section 2. Subsequently in Section 4 we will show how this specification can be correctly instrumented via program rewriting into a target language, hence we refer to the language introduced in this Section as our source language.

### 3.1. Source Language

Our source language model is essentially Featherweight Java (FJ) [7] with minor extensions including base types and an abstract notion of *library* methods for base types. The latter is important for an adequate consideration of taint propagation (e.g. on strings) in our model. FJ is a functional core calculus that includes class hierarchy definitions, subtyping, dynamic dispatch, and other basic features of Java. An FJ program is an expression $\mathsf{e}$ which is executed given a static *class table CT* which maintains class definitions. To describe program execution we will define a small step operational semantics relation on expressions $\mathsf{e}$ which we will take as synonymous with configurations as defined previously.

$$L ::= \texttt{class C extends C} \{\overline{\texttt{C}} \, \overline{\texttt{f}}; \, \texttt{K} \, \overline{\texttt{M}}\} \qquad\qquad \textit{classdefinitions}$$

$$K ::= \texttt{C}(\overline{\texttt{C}} \, \overline{\texttt{f}})\{\texttt{super}(\overline{\texttt{f}}); \, \texttt{this}.\overline{\texttt{f}} = \overline{\texttt{f}}; \} \qquad\qquad \textit{constructors}$$

$$M ::= \texttt{C m}(\overline{\texttt{C}} \, \overline{\texttt{x}})\{\texttt{return e}; \} \qquad\qquad \textit{methods}$$

$$e ::= \texttt{x} \mid \texttt{e.f} \mid \texttt{e.m}(\overline{\texttt{e}}) \mid \texttt{new C}(\overline{\texttt{e}}) \mid \texttt{if e then e else e} \mid \texttt{C.m(e)} \qquad \textit{expressions}$$

$$E ::= [\,] \mid \texttt{E.f} \mid \texttt{E.m}(\overline{\texttt{e}}) \mid \texttt{v.m}(\overline{\texttt{v}}, \texttt{E}, \overline{\texttt{e}}') \mid \texttt{new C}(\overline{\texttt{v}}, \texttt{E}, \overline{\texttt{e}}') \mid \texttt{if E then e else e} \mid \texttt{C.m(E)} \quad \textit{evaluation contexts}$$

Fig. 1. FJ Syntax

### 3.1.1. Syntax

The syntax of FJ is defined in Figure 1. We let $\texttt{A}, \texttt{B}, \texttt{C}, \texttt{D}$ range over class names, $\texttt{x}$ range over variables, $\texttt{f}$ range over field names, and $\texttt{m}$ range over method names. *Values*, denoted $\texttt{v}$ or $\texttt{u}$, are objects, i.e. expressions of the form $\texttt{new C}(\texttt{v}_1, \ldots, \texttt{v}_n)$. We assume given an $\texttt{Object}$ value that has no fields or methods. In addition to the standard expressions of FJ, we introduce a new form $\texttt{C.m(e)}$. This form is used to identify the method $\texttt{C.m}$ associated with a current evaluation context (aka the "activation frame"). This does not really change the semantics, but is a useful feature for our specification of sanitizer endorsement since return values from sanitizers need to be endorsed– see the Invoke and Return rules in the operational semantics below for its usage.

Conditional expressions are an important feature of the language for this presentation, since they are a control flow operation that should not be considered in a direct flow analysis. We assume that in any program setting true and false values, denote **T** and **F**, will be specified. When we consider base values and library methods below in Section 3.1.6, we will define a particular boolean value that we will use in this presentation.

For brevity in this syntax, we use vector notations. Specifically we write $\overline{\texttt{f}}$ to denote the sequence $\texttt{f}_1, \ldots, \texttt{f}_n$, similarly for $\overline{\texttt{C}}, \overline{\texttt{m}}, \overline{\texttt{x}}, \overline{\texttt{e}}$, etc., and we write $\overline{\texttt{M}}$ as shorthand for $\texttt{M}_1 \cdots \texttt{M}_n$. We write the empty sequence as $\varnothing$, we use a comma as a sequence concatenation operator. If and only if $\texttt{m}$ is one of the names in $\overline{\texttt{m}}$, we write $\texttt{m} \in \overline{\texttt{m}}$. Vector notation is also used to abbreviate sequences of declarations; we let $\overline{\texttt{C}} \, \overline{\texttt{f}}$ and $\overline{\texttt{C}} \, \overline{\texttt{f}};$ denote $\texttt{C}_1 \, \texttt{f}_1, \ldots, \texttt{C}_n \, \texttt{f}_n$ and $\texttt{C}_1 \, \texttt{f}_1; \ldots; \texttt{C}_n \, \texttt{f}_n;$ respectively. The notation $\texttt{this}.\overline{\texttt{f}} = \overline{\texttt{f}};$ abbreviates $\texttt{this}.\texttt{f}_1 = \texttt{f}_1; \ldots; \texttt{this}.\texttt{f}_n = \texttt{f}_n;$. Sequences of names and declarations are assumed to contain no duplicate names.

### 3.1.2. The class table and field and method body lookup

The class table *CT* maintains class definitions. The manner in which we look up field and method definitions implements inheritance and override, which allows fields and methods to be redefined in subclasses. Given a class table *CT*, the definitions of $mbody_{CT}(\texttt{C}, \texttt{m})$ and $fields_{CT}(\texttt{C})$ are given in Figure 2.

### 3.1.3. Method type lookup

Just as we've defined a function for looking up method bodies in the class table, we also define a function $mtype_{CT}(\texttt{C}, \texttt{m})$ that will look up types of a method $\texttt{C.m}$ in a class table in Figure 2. Although we omit FJ type analysis from this presentation, method type lookup will be useful for taint analysis instrumentation (Definition 4.1).

### 3.1.4. Operational semantics

Now, we can define the operational semantics of FJ. The reduction relation is binary, of the form $\kappa \to \kappa'$, and is defined via the inference rules in Figure 3.

$$fields_{CT}(\texttt{Object}) = \varnothing \qquad \frac{CT(\texttt{C}) = \texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad fields_{CT}(\texttt{D}) = \overline{\texttt{D}}\ \overline{\texttt{g}}}{fields_{CT}(\texttt{C}) = \overline{\texttt{D}}\ \overline{\texttt{g}}, \overline{\texttt{C}}\ \overline{\texttt{f}}}$$

$$\frac{CT(\texttt{C}) = \texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{B m}(\overline{\texttt{B}}\ \overline{\texttt{x}})\{\texttt{return e;}\} \in \overline{\texttt{M}}}{mbody_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \texttt{e}}$$

$$\frac{CT(\texttt{C}) = \texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{m} \notin \overline{\texttt{M}}}{mbody_{CT}(\texttt{m},\texttt{C}) = mbody_{CT}(\texttt{m},\texttt{D})}$$

$$\frac{\texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{B m}(\overline{\texttt{B}}\ \overline{\texttt{x}})\{\texttt{return e;}\} \in \overline{\texttt{M}}}{mtype_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{B}} \to \texttt{B}} \qquad \frac{\texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{m} \notin \overline{\texttt{M}}}{mtype_{CT}(\texttt{m},\texttt{C}) = mtype_{CT}(\texttt{m},\texttt{D})}$$

Fig. 2. Object Field, Method Body, and Method Type Lookup

Context
$$\frac{\texttt{e} \to \texttt{e}'}{\texttt{E[e]} \to \texttt{E[e}']}$$

Field
$$\frac{fields_{CT}(\texttt{C}) = \overline{\texttt{C}}\ \overline{\texttt{f}} \qquad \texttt{f}_i \in \overline{\texttt{f}}}{\texttt{new C}(\overline{\texttt{v}}).\texttt{f}_i \to \texttt{v}_i}$$

Invoke
$$\frac{mbody_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \texttt{e}}{\texttt{new C}(\overline{\texttt{v}}).\texttt{m}(\overline{\texttt{u}}) \to \texttt{C}.\texttt{m}(\texttt{e}[\texttt{new C}(\overline{\texttt{v}})/\texttt{this}][\overline{\texttt{u}}/\overline{\texttt{x}}])}$$

IfT
$$\texttt{if } \mathbf{T} \texttt{ then } \texttt{e}_1 \texttt{ else } \texttt{e}_2 \to \texttt{e}_1$$

IfF
$$\texttt{if } \mathbf{F} \texttt{ then } \texttt{e}_1 \texttt{ else } \texttt{e}_2 \to \texttt{e}_2$$

Return
$$\texttt{C}.\texttt{m}(\texttt{v}) \to \texttt{v}$$

Fig. 3. Operational Semantics for FJ

The definition of $\to$ assumes given a class table *CT* which is typically clear from context, but we will write $CT \vdash \kappa \to \kappa'$ to disambiguate class tables used in reductions when necessary. The definition also assumes that boolean values $\mathbf{T}$ and $\mathbf{F}$ are specified. We use $\to^*$ to denote the reflexive, transitive closure of $\to$, and we use $\to^n$ to denote an *n*-step reduction. We will also use the notion of an *execution trace* $\tau$ to range over sequences of configurations $\kappa_0 \ldots \kappa_n$ where $\kappa_i \to \kappa_{i+1}$ for all $0 \leqslant i < n$. Note that an execution trace $\tau$ may represent the partial execution of a program, i.e. the trace $\tau$ may be extended with additional configurations as the program continues execution. In general we will write $CT \vdash_\to \tau$ to disambiguate the class table *CT* and reduction relation $\to$ used for a trace $\tau$ when it is not clear from context.

### 3.1.5. Top-Level Programs

We define *top-level programs* $\mathfrak{p}(\mathbf{a})$ as programs of the form:

$$\texttt{new TopLevel().main}(\mathbf{a})$$

where $\mathbf{a}$ is a primitive object $\texttt{new C}(\overline{v})$. We assume that all class tables *CT* include an entry point $\texttt{TopLevel.main}$ with formal parameter $\texttt{attack}$, where $\texttt{TopLevel}$ objects have no fields. We write $\mathfrak{p}(\mathbf{a}) \Downarrow \tau$ iff trace $\tau$ begins with the configuration $\mathfrak{p}(\mathbf{a})$.

### 3.1.6. Library Methods

In order to study dynamic integrity taint analysis in FJ, we extend the semantics for library methods that allow specification of operations on base values (such as strings and integers). Consideration of these features is important for a thorough modeling of Phosphor-style taint analysis, and important related issues such as string- vs. character-based taint [23] which have not been considered in previous formal work on taint analysis [1]. Since static analysis is not a topic of this paper, for brevity we omit the standard FJ type analysis which is described in [7].

The abstract calculus described above is not particularly interesting with respect to direct information flow and integrity propagation, especially since method dispatch and conditional expressions are control flows that are discounted in direct data flow. More interesting is the manner in which taint propagates through base values and library operations, since direct flows propagate through some of these methods. Also, for run-time efficiency and ease of coding some Java taint analysis tools treat even complex library methods as "black boxes" that are instrumented at the top level for efficiency [24], rather than relying on instrumentation of lower-level operations.

Note that treating library methods as "black boxes" introduces a potential for over- and under-tainting– for example in some systems all string library methods that return strings are instrumented to return tainted results if any of the arguments are tainted, regardless of any direct flow from the argument to result [24]. Clearly this strategy introduces a potential for over-taint. Other systems do not propagate taint from strings to their component characters when decomposed [23], which is an example of under-taint. Part of our goal here is to develop an adequate language model to consider these approaches.

We therefore extend our basic definitions to accommodate base values and their manipulation. Let a *primitive field* be a field containing a base value. We call a *base type* any class with primitive fields only, and a *library method* is any method that operates on base type objects, defined in a primitive class. We expect primitive objects to be object wrappers for primitive values (e.g., $\texttt{Int}(5)$ wrapping primitive value $5$), and library methods to be object-oriented wrappers over primitive operations (e.g., $\texttt{Int plus(Int)}$ wrapping primitive operation $+$), allowing the latter's embedding in FJ. As a sanity condition we only allow library methods to select primitive fields or perform primitive operations. Let *LibMeths* be the set of library method names paired with their corresponding base classes in *BaseTypes*.

We posit a special set of field names *PrimField* that access primitive values ranged over by $v$ that may occur in objects, and a set of operations ranged over by $Op$ that operate on primitive values. We require that special field name selections only occur as arguments to $Op$, which can easily be enforced in practice by a static analysis. Similarly, primitive values $v$ may only occur in special object fields and be manipulated there by any $Op$.

$$\texttt{f}^* \in \textit{PrimField}$$
$$e ::= v \mid \texttt{e.f}^*$$
$$\texttt{e} ::= \cdots \mid Op(\bar{e})$$
$$\texttt{v} ::= \texttt{new C}(\overline{\texttt{v}}) \mid v$$
$$\texttt{E} ::= \cdots \mid Op(\bar{v}, \texttt{E}, \bar{e})$$

For library methods we require that the body of any library method be of the form where $\texttt{C}$ is a primitive class:

$$\texttt{return new C}(\overline{\texttt{e}}_1, \ldots, \overline{\texttt{e}}_n)$$

We define the meaning of operations *Op* via an "immediate" big-step semantic relation $\approx$ where the rhs of the relation is required to be a primitive value, and we identify expressions up to $\approx$. For example, to define a library method for integer addition, where `Int` objects contain a primitive numeric `val`, field we would define a $+$ operation as follows:

$$+(n_1, n_2) \approx n_1 + n_2$$

Then we can add to the definition of `Int` in *CT* a method `Plus` to support arithmetic in programs:

```
Int plus(Int x) { return(new Int(+(this.val, x.val))); }
```

Similarly, to define string concatenation, we define a concatenation operation @ on primitive strings:

$$@(s_1, s_2) \approx s_1 s_2$$

and we extend the definition of `String` in *CT* with the following method, where we assume all `String` objects maintain their primitive representation in a `val` field:

```
String concat(String x)
        { return(new String(@(this.val, x.val))); }
```

A boolean class `Bool` can be defined on the basis of constants *true* and *false* and standard boolean connectives– we will subsequently use this encoding for values **T** and **F** and conditional guards:

$$b \in \{true, false\} \qquad \mathbf{T} \triangleq \texttt{new Bool}(true) \qquad \mathbf{F} \triangleq \texttt{new Bool}(false) \qquad \wedge(b_1, b_2) \approx b_1 \wedge b_2$$

$$\vee(b_1, b_2) \approx b_1 \vee b_2 \qquad \neg(b) \approx \neg b$$

These boolean values can represent the results of base object comparison operators such as a string equality test:

$$eq(s_1, s_2) \approx b = \begin{cases} true \text{ if } s_1 = s_2 \\ false \text{ otherwise} \end{cases}$$

```
String eq(String x)
        { return(new Bool(eq(this.val, x.val))); }
```

### 3.2. In-Depth Integrity Analysis Specified Logically

In this section, we demonstrate how in-depth integrity taint analysis for FJ can be expressed as a single uniform policy separate from code. To accomplish this we interpret program traces as information represented by a logical fact base in the style of Datalog. We then define a predicate called Shadow that inductively constructs a "shadow" of configurations reflecting the taint of values.

$$toFOL(\mathtt{v}, n) = \{\text{Value}(n, \mathtt{v})\}$$

$$toFOL(\mathtt{E}[\mathtt{new}\ \mathtt{C}(\overline{\mathtt{v}}).\mathtt{f}], n) = \{\text{GetField}(n, \mathtt{new}\ \mathtt{C}(\overline{\mathtt{v}}), \mathtt{f}), \text{Context}(n, \mathtt{E})\}$$

$$toFOL(\mathtt{E}[\mathtt{new}\ \mathtt{C}(\overline{\mathtt{v}}).\mathtt{m}(\overline{\mathtt{u}})], n) = \{\text{Call}(n, \mathtt{C}, \overline{\mathtt{v}}, \mathtt{m}, \overline{\mathtt{u}}), \text{Context}(n, \mathtt{E})\}$$

$$toFOL(\mathtt{E}[\mathtt{C}.\mathtt{m}(\mathtt{v})], n) = \{\text{ReturnValue}(n, \mathtt{C}, \mathtt{m}, \mathtt{v}), \text{Context}(n, \mathtt{E})\}$$

$$toFOL(\mathtt{E}[Op(\overline{v})], n) = \{\text{PrimCall}(n, Op, \overline{v}), \text{Context}(n, \mathtt{E})\}$$

$$toFOL(\mathtt{E}[\mathtt{if}\ \mathbf{T}\ \mathtt{then}\ \mathtt{e}_1\ \mathtt{else}\ \mathtt{e}_2], n) = \{\text{IfT}(n, \mathtt{e}_1, \mathtt{e}_2), \text{Context}(n, \mathtt{E})\}.$$

$$toFOL(\mathtt{E}[\mathtt{if}\ \mathbf{F}\ \mathtt{then}\ \mathtt{e}_1\ \mathtt{else}\ \mathtt{e}_2], n) = \{\text{IfF}(n, \mathtt{e}_1, \mathtt{e}_2), \text{Context}(n, \mathtt{E})\}.$$

Fig. 4. Interpreting Expressions as Formulas via $toFOL(\cdot)$.

Java-based taint analyses naturally tend to be object-based, i.e. low-integrity values are objects conceptually, and objects have an assigned taint level in the implementation. The types of tainted objects varies depending on the analysis, but most emphasize taint of base values. We will likewise focus on taint of base values, though we will support taint labeling of *all* objects. This is partly to generalize the representation, but also for formal convenience– In our logical specification of taint analysis, a shadow expression has a syntactic structure that matches the configuration expression, and associates integrity levels (including "high" ∘ and "low" •) with particular objects via shape conformance.

**Example 3.1.** *Suppose a method* m *of an untainted* C *object with no fields is invoked on a pair of tainted* $s_1$ *and untainted* $s_2$ *strings:*

$$\mathtt{new}\ \mathtt{C}().\mathtt{m}(\mathtt{new}\ \mathtt{String}(s_1), \mathtt{new}\ \mathtt{String}(s_2))$$

*Its proper shadow is:*

$$\mathtt{shadow}\ \mathtt{C}(\circ).\mathtt{m}(\mathtt{shadow}\ \mathtt{String}(\bullet), \mathtt{shadow}\ \mathtt{String}(\circ)).$$

On the basis of shadow expressions that correctly track integrity, we can logically specify prospective taint analysis as a property of shadowed trace information, and retrospective taint analysis as a function of shadowed trace information. An extended example of a shadowed trace is presented in Section 3.2.4.

*3.2.1. Taint Tracking as a Logical Trace Property*

In order to specify taint tracking, we define the mapping $toFOL(\cdot)$ that shows how we concretely model execution traces in FOL. We develop $toFOL(\cdot)$ that interprets FJ traces as sets of logical facts (a fact base). Intuitively, in the interpretation each configuration is represented by a Context predicate representing the evaluation context, and a predicate representing the redex (e.g. Call). Each of these predicates has an initial natural number argument denoting a "timestamp" that orders configurations in a trace.

**Definition 3.1.** *We define* $toFOL(\cdot)$ *as a mapping on traces and configurations:*

$$toFOL(\tau) = \bigcup_{\sigma \in \mathbf{prefix}(\tau)} toFOL(\sigma)$$

*such that $toFOL(\sigma) = \bigcup_i toFOL(\kappa_i, i)$ for $\sigma = \kappa_1 \cdots \kappa_k$. We define $toFOL(\kappa, n)$ in Figure 4.*

*Integrity Identifiers.*  We introduce an integrity identifier $t$ that denotes the integrity level associated with objects. To support a notion of "partial endorsement" for partially trusted sanitizers, we define three taint labels, to denote high integrity ($\circ$), low integrity ($\bullet$), and uncertain integrity ($\odot$). We refer to these levels as tainted, untainted, and maybe tainted, respectively.

$$t ::= \circ \mid \odot \mid \bullet$$

We specify an ordering $\leqslant$ on these labels denoting their integrity relation:

$$\bullet \leqslant \odot \leqslant \circ$$

For simplicity in this presentation we will assume that all *Sanitizers* are partially trusted and cannot raise the integrity of a tainted or maybe tainted object beyond maybe tainted. It would be possible to include both trusted and untrusted sanitizers without changing the formalism.

We posit the usual meet $\wedge$ and join $\vee$ operations on taint lattice elements, and introduce logical predicates meet and join such that $meet(t_1 \wedge t_2, t_1, t_2)$ and $join(t_1 \vee t_2, t_1, t_2)$ hold.

### 3.2.2. Shadow Traces, Taint Propagation, and Sanitization

Shadow traces reflect taint information of objects as they are passed around programs. Shadow traces are comprised of shadow expressions and contexts which are terms in the logic with the following syntax. Note the structural conformance with closed e and E, but with primitive values replaced with a single dummy value $\delta$ that is omitted for brevity in examples, but is necessary to maintain proper arity for field selection. Shadow expressions most importantly assign integrity identifiers $t$ to objects:

$$sv ::= \texttt{shadow } C(t, \overline{sv}) \mid \delta$$
$$se ::= sv \mid se.\texttt{f} \mid se.\texttt{m}(\overline{se}) \mid \texttt{shadow } C(t, \overline{se}) \mid C.\texttt{m}(se) \mid Op(\overline{se}) \mid \texttt{if } se \texttt{ then } se \texttt{ else } se$$
$$SE ::= [\,] \mid SE.\texttt{f} \mid SE.\texttt{m}(\overline{se}) \mid sv.\texttt{m}(\overline{sv}, SE, \overline{se}') \mid \texttt{shadow } C(t, \overline{sv}, SE, \overline{se}') \mid C.\texttt{m}(SE) \mid$$
$$\qquad Op(\overline{sv}, SE, \overline{se}) \mid \texttt{if } SE \texttt{ then } se \texttt{ else } se$$

The shadowing specification requires that shadow expressions evolve in a shape-conformant way with the original configuration. To this end, we define a metatheoretic function for shadow method bodies, *smbody*, that imposes untainted tags on all method bodies, defined a priori, and removes primitive values.

**Definition 3.2.** *Shadow method bodies are defined by the function smbody.*

$$smbody_{CT}(\texttt{m}, C) = \overline{\texttt{x}}.srewrite(\texttt{e}),$$

*where $mbody_{CT}(\texttt{m}, C) = \overline{\texttt{x}}.\texttt{e}$ and the shadow rewriting function, srewrite, is defined as follows, where $srewrite(\overline{\texttt{e}})$ denotes a mapping of srewrite over the vector $\overline{\texttt{e}}$:*

$$srewrite(\texttt{x}) = \texttt{x}$$
$$srewrite(\texttt{new } C(\overline{\texttt{e}})) = \texttt{shadow } C(\circ, srewrite(\overline{\texttt{e}}))$$
$$srewrite(\texttt{e.f}) = srewrite(\texttt{e}).\texttt{f}$$

$\text{match}(sv, [\,], sv).$

$\text{match}(\texttt{shadow C}(t, \overline{sv}).\texttt{f}_i, [\,], \texttt{shadow C}(t, \overline{sv}).\texttt{f}_i).$

$\text{match}(\texttt{shadow C}(t, \overline{sv}).\texttt{m}(\overline{su}), [\,], \texttt{shadow C}(t, \overline{sv}).\texttt{m}(\overline{su})).$

$\text{match}(\texttt{C}.\texttt{m}(sv), [\,], \texttt{C}.\texttt{m}(sv)).$

$\text{match}(\texttt{if } sv \texttt{ then } se_1 \texttt{ else } se_2, [\,], \texttt{if } sv \texttt{ then } se_1 \texttt{ else } se_2).$

$\text{match}(se, SE, se') \implies \text{match}(se.\texttt{f}, SE.\texttt{f}, se').$

$\text{match}(se, SE, se') \implies \text{match}(se.\texttt{m}(\overline{se}), SE.\texttt{m}(\overline{se}), se').$

$\text{match}(se, SE, se') \implies \text{match}(sv.\texttt{m}(\overline{sv}, se, \overline{se}), sv.\texttt{m}(\overline{sv}, SE, \overline{se}), se').$

$\text{match}(se, SE, se') \implies \text{match}(\texttt{shadow C}(t, \overline{sv}, se, \overline{se}), \texttt{shadow C}(t, \overline{sv}, SE, \overline{se}), se').$

$\text{match}(se, SE, se') \implies \text{match}(\texttt{C}.\texttt{m}(se), \texttt{C}.\texttt{m}(SE), se').$

$\text{match}(se, SE, se') \implies \text{match}(Op(\overline{sv}, se, \overline{se}), Op(\overline{sv}, SE, \overline{se}), se').$

$\text{match}(se, SE, se') \implies \text{match}(\texttt{if } se \texttt{ then } se_1 \texttt{ else } se_2, \texttt{if } SE \texttt{ then } se_1 \texttt{ else } se_2, se')$

Fig. 5. match Predicate Definition.

$srewrite(\texttt{e}.\texttt{m}(\overline{\texttt{e}}')) = srewrite(\texttt{e}).\texttt{m}(srewrite(\overline{\texttt{e}}'))$

$srewrite(\texttt{C}.\texttt{m}(\texttt{e})) = \texttt{C}.\texttt{m}(srewrite(\texttt{e}))$

$srewrite(Op(\overline{\texttt{e}})) = Op(srewrite(\overline{\texttt{e}}))$

$srewrite(\texttt{if } \texttt{e}_1 \texttt{ then } \texttt{e}_2 \texttt{ else } \texttt{e}_3) = \texttt{if } srewrite(\texttt{e}_1) \texttt{ then } srewrite(\texttt{e}_2) \texttt{ else } srewrite(\texttt{e}_3)$

$srewrite(v) = \delta$

We use match as a predicate which matches a shadow expression *se*, to a shadow context *SE* and a shadow expression *se'* where *se'* is the part of the shadow in the hole. The definition of match is given in Figure 5.

Next, in Figure 6, we define a predicate $\text{Shadow}(n, se)$ where *se* is the relevant shadow expression at execution step *n*, establishing an ordering for the shadow trace. Shadow has as its precondition a "current" shadow expression, and as its postcondition the shadow expression for the next step of evaluation (with the exception of the rule for shadowing *Op*s on primitive values which reflects the "immediate" valuation due to the definition of ≈– note the timestamp is not incremented in the postcondition in that case). We set the shadow of the initial configuration at timestamp 1, and then Shadow inductively shadows the full trace. Shadow is defined by case analysis on the structure of shadow expression in the hole. The shadow expression in the hole and the shadow evaluation context are derived from match predicate definition.[3]

With respect to control flow, the most notable rules of Shadow are those governing conditional branching, which ignore the taint of the guard, and method dispatch, which ignore the taint of the object associated with the dispatched method. Since we focus on base value taint, method dispatch is essentially a

---

[3]Some notational liberties are taken in Figure 6 regarding expression and context substitutions, which are defined using predicates elided for brevity.

non-issue, however conditional branching is directly dependent on base values so ignoring the taint of the guard explicitly ignores indirect data flow.

*Taint Propagation and Endorsement..* The propagation of taint in the model described in Section 1.2 is embedded in the definition of Shadow, in particular we assume a set of *Sanitizers*. For elements of *Sanitizers*, if input is tainted then the result is considered to be only partially endorsed (maybe tainted). For library methods, taint is propagated given a user-defined predicate $\text{Prop}(t, \iota)$ where $\iota$ is a compound term of the form $\texttt{C.m}(\bar{t})$ with $\bar{t}$ the given integrity of $\texttt{this}$ followed by the integrity of the arguments to method $\texttt{C.m}$, and $t$ is the integrity of the result. For example, one could define:

$$\text{meet}(t, t_1, t_2) \Rightarrow \text{Prop}(t, \texttt{String.concat}(t_1, t_2)) \qquad \text{meet}(t, t_1, t_2) \Rightarrow \text{Prop}(t, \texttt{String.eq}(t_1, t_2))$$

Later in Section 5.2.1 we will discuss formal semantic conditions on library methods that ensure sound taint propagation.

### 3.2.3. In-Depth Integrity Taint Analysis Policies

Now we can logically specify an in-depth policy for integrity taint analysis, as proposed originally in Section 1.2. In particular we assume a set *Sanitizers* and a set *SSOs*. Since objects may inherit a sanitizer or sso from a superclass, we require that *Sanitizers* and *SSOs* are closed under inheritance as a sanity condition, as follows:

$$\frac{CT(\texttt{C}) = \texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{m} \notin \overline{\texttt{M}} \qquad \texttt{D.m} \in \textit{SSOs}}{\texttt{C.m} \in \textit{SSOs}}$$

$$\frac{CT(\texttt{C}) = \texttt{class C extends D} \{\overline{\texttt{C}}\ \overline{\texttt{f}};\ \texttt{K}\ \overline{\texttt{M}}\} \qquad \texttt{m} \notin \overline{\texttt{M}} \qquad \texttt{D.m} \in \textit{Sanitizers}}{\texttt{C.m} \in \textit{Sanitizers}}$$

The in-depth policy has both prospective and retrospective component– the former is defined as a safety property [42], while the latter is defined as a logging specification. The prospective component of the policy must identify traces where a tainted value is passed to a secure method. To this end, in Figure 7 we define the predicate BAD which identifies traces that should be rejected as unsafe– a bad trace is any in which an sso is executed with a tainted argument. The retrospective component specifies that data of questionable integrity that is passed to a secure method should be logged. The relevant logging specification is specified in terms of a predicate MaybeBAD also defined in Figure 7.

**Definition 3.3.** *Let X be the set of rules in Figures 5, 6, and 7 and the set of user-defined rules for* Prop. *The prospective integrity taint analysis policy is defined as the set of traces that are free from, or at most end in,* BAD *configurations. This latter condition is not necessary for the specification, and may seem extraneous, but it is mainly in place to allow a clean proof correspondence with the implementation (as detailed in Section 4.2), since taint checks will be placed to execute immediately* after *sso invocation in the implementation.*

$$\text{SP}_{\text{taint}} = \{\tau\kappa \mid (\lfloor\tau\rfloor \otimes C(X))^{\Rightarrow\{\text{BAD}\}} = C(\varnothing)\}.$$

*The retrospective integrity taint analysis policy is the following logging specification– in this definition again to simplify semantic correspondence with instrumentation, we specify that arguments should be*

Shadow(1, shadow TopLevel(o).main(shadow C($\bullet, \overline{\delta}$))).

Shadow($n, se$) $\land$ match($se, SE, sv$.m($\overline{sv}'$)) $\land$ C.m $\notin$ *LibMeths* $\land$ *smbody*$_{CT}$(m, C) $= \overline{\text{x}}.se'$ $\implies$

    Shadow($n+1, SE$[C.m($se'[\overline{sv}'/\overline{x}][sv/\texttt{this}]$)]).

Shadow($n, se$) $\land$ match($se, SE$, shadow C($t_0, \overline{sv}$).m($\overline{\text{shadow C}(t, \overline{sv})}$)) $\land$ C.m $\in$ *LibMeths* $\land$ *smbody*$_{CT}$(m, C) $= \overline{\text{x}}.$shadow D(o, $\overline{se}$)$\land$

    Prop($t$, C.m($t_0, \overline{t}$)) $\implies$ Shadow($n+1, SE$[C.m(shadow D($t, \overline{se}$)[shadow C($t_0, \overline{sv}$)/this][$\overline{\text{shadow C}(t, \overline{sv})}/\overline{\text{x}}$])]).

Shadow($n, se$) $\land$ match($se, SE$, shadow C($t, \overline{sv}$).f$_\texttt{i}$) $\implies$ Shadow($n+1, SE$[$sv_i$]).

Shadow($n, se$) $\land$ match($se, SE, Op(\overline{\delta})$) $\implies$ Shadow($n, SE[\delta]$).

Shadow($n, se$) $\land$ match($se, SE$, C.m(shadow D($t, \overline{sv}$))) $\land$ C.m $\in$ *Sanitizers* $\implies$ Shadow($n+1, SE$[shadow D($t \lor \odot, \overline{sv}$)]).

Shadow($n, se$) $\land$ match($se, SE$, C.m($sv$)) $\land$ C.m $\notin$ *Sanitizers* $\implies$ Shadow($n+1, SE[sv]$).

Shadow($n, se$) $\land$ match($se, SE$, if $sv$ then $se_1$ else $se_2$) $\land$ IfT($n, \texttt{e}_1, \texttt{e}_2$) $\implies$ Shadow($n+1, SE[se_1]$).

Shadow($n, se$) $\land$ match($se, SE$, if $sv$ then $se_1$ else $se_2$) $\land$ IfF($n, \texttt{e}_1, \texttt{e}_2$) $\implies$ Shadow($n+1, SE[se_2]$)

Fig. 6. Shadow Predicate Definition.

*logged one step after invocation of an sso:*

$$LS_{\text{taint}} = \lambda\tau\kappa.(\lfloor\tau\rfloor \otimes C(X))^{\Rightarrow\{\text{MaybeBAD}\}}$$

We immediately observe that SP$_{\text{taint}}$ is a safety property:

**Lemma 3.1.** SP$_{\text{taint}}$ *is a safety property.*

**Proof.** Let $\tau\kappa \notin$ SP$_{\text{taint}}$. Then, $(\lfloor\tau\rfloor \otimes C(X))^{\Rightarrow\{\text{BAD}\}} \neq C(\varnothing)$. This implies that there exists some $n$ such that BAD($n$) $\in (\lfloor\tau\rfloor \otimes C(X))^{\Rightarrow\{\text{BAD}\}}$. Let $\tau[\cdots n]$ denote the finite prefix of $\tau$ up to timestamp $n$. By Definition 3.3 BAD only refers to events that precede step $n$, so it follows that $\lfloor\tau\rfloor \otimes C(X) \vdash$ BAD($n$) iff $\lfloor\tau[\cdots n]\rfloor \otimes C(X) \vdash$ BAD($n$), i.e. $\tau \notin$ SP$_{\text{taint}}$ iff $\tau[\cdots n] \notin$ SP$_{\text{taint}}$ for finite $n$, hence SP$_{\text{taint}}$ is a safety property [42]. $\square$

Finally we define a program as being safe iff it does not produce a bad trace.

**Definition 3.4.** *We call a program* $\mathfrak{p}(\mathbf{a})$ safe *iff for all $\tau$ it is the case that* $\mathfrak{p}(\mathbf{a}) \Downarrow \tau$ *implies* $\tau \in$ SP$_{\text{taint}}$. *We call the program* unsafe *iff there exists some trace $\tau$ such that* $\mathfrak{p}(\mathbf{a}) \Downarrow \tau$ *and* $\tau \notin$ SP$_{\text{taint}}$.

### 3.2.4. Extended Example: Reduction and Shadowing

To illustrate the major points of our construction for source program traces and their shadows, we consider an example of program that contains an *sso* call on a string that has been constructed from a sanitized low integrity input.

**Example 3.2.** *Assume that sanitizer and sso methods* Sec.sanitize *and* Sec.secureMeth *are identity functions for the sake of brevity, i.e.:*

$$mbody_{CT}(\text{Sec}, \texttt{sanitize}) = \text{x}, \text{x} \qquad mbody_{CT}(\text{Sec}, \texttt{secureMeth}) = \text{x}, \text{x}$$

$$\mathrm{match}(se, SE, \texttt{shadow C}(t, \overline{sv}).\texttt{m}(\texttt{shadow D}(t', \overline{sv}')) \wedge \mathrm{Shadow}(n, se) \wedge \mathrm{Call}(n, \texttt{C}, \overline{v}, \texttt{m}, \texttt{u}) \wedge \texttt{C.m} \in SSOs \implies \mathrm{SsoTaint}(n, t', \texttt{u}).$$

$$\mathrm{SsoTaint}(n, \bullet, \texttt{u}) \implies \mathrm{BAD}(n). \qquad\qquad \mathrm{SsoTaint}(n, t, \texttt{u}) \wedge t \leqslant \odot \implies \mathrm{MaybeBAD}(\texttt{u}).$$

Fig. 7. Predicates for Specifying Prospective and Retrospective Properties

$$\mathfrak{p}(\texttt{new String}(''\texttt{hello}''))$$
$$\rightarrow^5$$
$$\texttt{TopLevel.main}(\texttt{new Sec}().\texttt{secureMeth}(\texttt{new Sec}().\texttt{sanitize}(\texttt{new String}(''\texttt{hello world}''))))$$
$$\rightarrow^2$$
$$\texttt{TopLevel.main}(\texttt{new Sec}().\texttt{secureMeth}(\texttt{new String}(''\texttt{hello world}'')))$$
$$\rightarrow^2$$
$$\texttt{TopLevel.main}(\texttt{new String}(''\texttt{hello world}''))$$
$$\rightarrow$$
$$\texttt{new String}(''\texttt{hello world}'')$$

Fig. 8. Example 3.2: Source Trace.

$$\mathrm{Shadow}\big(1, \texttt{shadow TopLevel}(\circ).\texttt{main}(\texttt{shadow String}(\bullet, \delta))\big)$$
$$\mathrm{Shadow}\big(5, \texttt{TopLevel.main}(\texttt{shadow Sec}(\circ).\texttt{secureMeth}(\texttt{shadow Sec}(\circ).\texttt{sanitize}(\texttt{shadow String}(\bullet, \delta))))\big)$$
$$\mathrm{Shadow}\big(7, \texttt{TopLevel.main}(\texttt{shadow Sec}(\circ).\texttt{secureMeth}(\texttt{shadow String}(\odot, \delta)))\big)$$
$$\mathrm{Shadow}\big(9, \texttt{TopLevel.main}(\texttt{shadow String}(\odot, \delta))\big)$$
$$\mathrm{Shadow}\big(10, \texttt{shadow String}(\odot, \delta)\big)$$

Fig. 9. Example 3.2: Shadow Expressions.

*and let* $mbody_{CT}(\texttt{main}, \texttt{TopLevel})$ *be:*

$$\texttt{attack}, \texttt{new Sec}().\texttt{secureMeth}(\texttt{new Sec}().\texttt{sanitize}(\texttt{attack.concat}(\texttt{new String}(''\texttt{world}'')))).$$

*Assume also that an input string* $''\texttt{hello}''$ *is tainted with low integrity– Figure 8 depicts a source trace given the initial configuration:*

$$\texttt{new TopLevel}().\texttt{main}(\texttt{new String}(''\texttt{hello}''))$$

*with some reduction steps elided to highlight calls to* $\texttt{Sec.sanitize}$ *and* $\texttt{Sec.secureMeth}$. *In Figure 9 we show shadows of configurations highlighted (depicted) in the source trace. We note this trace is in* $\mathrm{SP_{taint}}$ *and hence is safe.*

## 4. Correct Instrumentation via Program Rewriting

Now we define an object-based dynamic integrity taint analysis in a more familiar operational style. Taint analysis instrumentation is added automatically by a program rewriting algorithm *Phos* that models the Phosphor rewriting algorithm, defined in Section 4.1. It adds taint label fields to all objects, and operations for appropriately propagating taint along direct flow paths. In addition to blocking behavior to enforce prospective checks, we incorporate logging instrumentation to support retrospective measures in the presence of partially trusted sanitization. We illustrate computation of instrumented code via an extended example in Section 4.1.4, which continues the (now running) example introduced in Section 3.2.4

In Section 4.2 we follow the methods developed in Section 2 and show that *Phos* is semantics preserving, and that instrumented code generates sound and complete audit logs with respect to the logging specification $LS_{taint}$ defined in Section 3.2.3. We will also show that instrumented code respects the safety property $SP_{taint}$ defined in the latter Section.

### 4.1. In-Depth Taint Analysis Instrumentation

The target language $FJ_{taint}$ of the rewriting algorithm *Phos* has the same syntax as FJ except we add taint labels $t$ as a form of primitive value $v$, the type of which we posit as `Taint`. For the semantics of taint values operations we define:

$$\vee(t_1, t_2) \approx t_1 \vee t_2 \qquad\qquad \wedge(t_1, t_2) \approx t_1 \wedge t_2$$

In addition we introduce a "check" operation $?$ such that $?t \approx t$ iff $t > \bullet$. We also add a convenient sequencing operation of the form $e;e$ to target language expressions, and evaluation contexts of the form $E;e$.

#### 4.1.1. The Phos Algorithm.

Now we define the program rewriting algorithm *Phos* as follows. It incorporates a rewriting function $\mu$ that assigns an untainted label to every object in an FJ source program. The class table is manipulated by *Phos* to specify a `taint` field for all objects, a `check` object method that blocks if the argument is tainted, and an `endorse` method for any object returned by a sanitizer.

**Definition 4.1.** *For any expression* $e$*, the expression* $\mu(e)$ *is syntactically equivalent to* $e$ *except every subexpression* `new C(`$\overline{e}$`)` *is replaced with* `new C(o,`$\overline{e}$`)`*. Given SSOs and Sanitizers, define:*

$$Phos(e, CT) = (\mu(e), Phos(CT))$$

*where* $Phos(CT)$ *is the smallest class table satisfying the axioms given in Figure 10. Furthermore, to correctly mark low integrity input as tainted, given class table* $CT$ *and top-level program* $\mathfrak{p}(\mathbf{a})$ *where* $\mathbf{a} = $ `new C(`$\overline{v}$`)` *we define:*

$$Phos(\mathfrak{p}(\mathbf{a})) = Phos(\mathfrak{p}, CT)(\text{new C}(\bullet, \overline{v}))$$

As discussed in Section 1, sanitization is typically taken to be "ideal" for integrity flow analyses, however in practice sanitization is imperfect, which creates an attack vector. To support retrospective

$$\mathit{fields}_{\mathit{Phos}(CT)}(\texttt{Object}) = \texttt{Taint taint} \qquad \mathit{mbody}_{\mathit{Phos}(CT)}(\texttt{check},\texttt{Object}) = \texttt{x}, \texttt{new Object}(\texttt{?x.taint}); \texttt{x}$$

$$\frac{\texttt{C.m} \in \mathit{Sanitizers} \qquad \mathit{mtype}_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{C}} \rightarrow \texttt{D} \qquad \mathit{fields}_{CT}(\texttt{D}) = \overline{\texttt{f}}}{\mathit{mbody}_{\mathit{Phos}(CT)}(\texttt{endorse},\texttt{D}) \quad = \quad \varnothing, \texttt{new D}(\vee(\odot,\texttt{this.taint}),\overline{\texttt{this.f}})}$$

$$\frac{\texttt{C.m} \in \mathit{SSOs} \qquad \mathit{mbody}_{CT}(\texttt{m},\texttt{C}) = \texttt{x}, \texttt{e}}{\mathit{mbody}_{\mathit{Phos}(CT)}(\texttt{m},\texttt{C}) = \texttt{x}, \texttt{this.log(x)}; \texttt{this.check(x)}; \mu(\texttt{e})} \qquad \frac{\texttt{C.m} \in \mathit{Sanitizers} \qquad \mathit{mbody}_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \texttt{e}}{\mathit{mbody}_{\mathit{Phos}(CT)}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \mu(\texttt{e}).\texttt{endorse()}}$$

$$\frac{\texttt{C.m} \notin \mathit{Sanitizers} \cup \mathit{SSOs} \qquad \mathit{mbody}_{CT}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \texttt{e}}{\mathit{mbody}_{\mathit{Phos}(CT)}(\texttt{m},\texttt{C}) = \overline{\texttt{x}}, \mu(\texttt{e})}$$

Fig. 10. Axioms for Rewriting Algorithm

measures specified in Definition 3.3, we define `endorse` so it takes object taint *t* to the join of *t* and $\odot$. The algorithm also adds a `log` method call to the beginning of *SSOs*, which will log objects that are maybe tainted or worse. The semantics of `log` are defined directly in the operational semantics of $\text{FJ}_{\text{taint}}$ below.

### 4.1.2. Taint Propagation of Library Methods

Another important element of taint analysis is instrumentation of library methods that propagate taint– the propagation must be made explicit to reflect the interference of arguments with results. The approach to this in taint analysis systems is often motivated by efficiency as much as correctness [24]. We assume that library methods are instrumented to propagate taint as intended (i.e. in accordance with the user defined predicate Prop).

Here is how addition, string concatenation, and equality test, can be modified to propagate taint. Note the taint of arguments will be propagated to results by taking the meet of argument taint, thus reflecting the degree of integrity corruption:

```
Int plus(Int x)
    { return(new Int
      (∧(this.taint, x.taint), +(this.val, x.val))); }
```
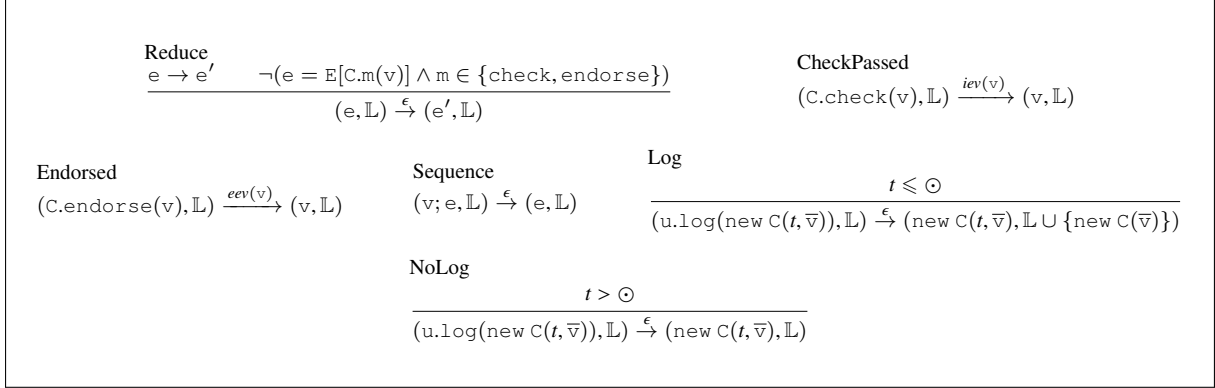
```
String concat(String x)
    { return(new String
      (∧(this.taint, x.taint), @(this.val, x.val))); }
```

```
String eq(String x)
    { return(new Bool
      (∧(this.taint, x.taint), eq(this.val, x.val))); }
```

### 4.1.3. Operational Semantics of FJ_taint

The operational semantics of $\text{FJ}_{\text{taint}}$ are defined in Figure 11. Configurations in FJ are of the form $(\texttt{e}, \mathbb{L})$ where reductions are defined in terms of a labeled transition relation $\xrightarrow{\alpha}$ on configurations, where

$$\frac{\text{Reduce}}{(\text{e}, \mathbb{L}) \xrightarrow{\epsilon} (\text{e}', \mathbb{L})} \qquad \neg(\text{e} = \text{E}[\text{C.m(v)}] \wedge \text{m} \in \{\text{check}, \text{endorse}\})$$

$$\text{CheckPassed} \qquad (\text{C.check(v)}, \mathbb{L}) \xrightarrow{iev(\text{v})} (\text{v}, \mathbb{L})$$

$$\text{Endorsed} \qquad (\text{C.endorse(v)}, \mathbb{L}) \xrightarrow{eev(\text{v})} (\text{v}, \mathbb{L}) \qquad \text{Sequence} \qquad (\text{v}; \text{e}, \mathbb{L}) \xrightarrow{\epsilon} (\text{e}, \mathbb{L})$$

$$\frac{\text{Log} \qquad t \leqslant \odot}{(\text{u.log(new C}(t, \overline{\text{v}})), \mathbb{L}) \xrightarrow{\epsilon} (\text{new C}(t, \overline{\text{v}}), \mathbb{L} \cup \{\text{new C}(\overline{\text{v}})\})}$$

$$\frac{\text{NoLog} \qquad t > \odot}{(\text{u.log(new C}(t, \overline{\text{v}})), \mathbb{L}) \xrightarrow{\epsilon} (\text{new C}(t, \overline{\text{v}}), \mathbb{L})}$$

Fig. 11. Operational Semantics of FJ$_{\text{taint}}$.

$\alpha$ is a possibly empty sequence $\epsilon$ of security events. These events are either *integrity* events $iev(\text{v})$ emitted when a check succeeds during evaluation as defined in the CheckPassed rule, or *endorsement* events $eev(\text{v})$, emitted when a value is endorsed as defined in the Endorsed rule.

Audit logs $\mathbb{L}$ are added to configurations to support the retrospective security via audit logging, and are defined as sets of objects (values). The log method is the only one that interacts with the log in any way, and its semantics are specified in the Log and NoLog rules, where possibly tainted values are logged, and untainted ones are not. Note that we strip taint tags from values for logging– this is mainly to simplify the correspondence with *LS*$_{\text{taint}}$ semantics for our technical development (where taint tags don't exist). We otherwise "inherit" the reduction semantics of FJ via the Reduce rule.

We write $\kappa_0 \xrightarrow{\alpha_0 \cdots \alpha_{n-1}}_n \kappa_n$ iff $\kappa_i \xrightarrow{\alpha_i} \kappa_{i+1}$ for all $0 \leqslant i < n$, and write $\kappa \xrightarrow{\alpha}_* \kappa'$ iff $\kappa \xrightarrow{\alpha}_n \kappa'$ for some $n$. We may omit transition labels in cases where they are empty ($\epsilon$) or not relevant to discussion, abusing notation $\rightarrow$, $\rightarrow^n$, and $\rightarrow^*$ as defined for FJ. We define traces as for FJ, and we write $\text{e} \Downarrow \tau$ iff $\tau$ begins with the configuration $(\text{e}, \emptyset)$.

### 4.1.4. Extended Example: Target Trace

Revisiting the example introduced in Section 3.2.4, we show execution of the rewritten program *Phos*($\mathfrak{p}(\mathbf{a})$) in Figure 12. By definition the rewritten top-level program is:

$$\text{new TopLevel(o).main(new String}(\bullet, ''\text{hello}''))$$

We note that additional reduction steps are necessary to evaluate instrumentation code in the target program, and that *eev* and *iev* events mark points during reduction when a value is endorsed and when it is checked.

### 4.2. Operational Properties of Phos

Now we can leverage machinery developed previously to demonstrate in-depth operational correctness of *Phos*, i.e. both prospective and retrospective operational correctness.

#### 4.2.1. Main Results

Recalling our definitions of semantics preservation, soundness, and completeness from Section 2, we state our main results as follows. These results tie together our relevant logging specification *LS*$_{\text{taint}}$ and safety property SP$_{\text{taint}}$ defined in Section 3.2.3. Preliminary proof details are presented in Sections 4.2.2
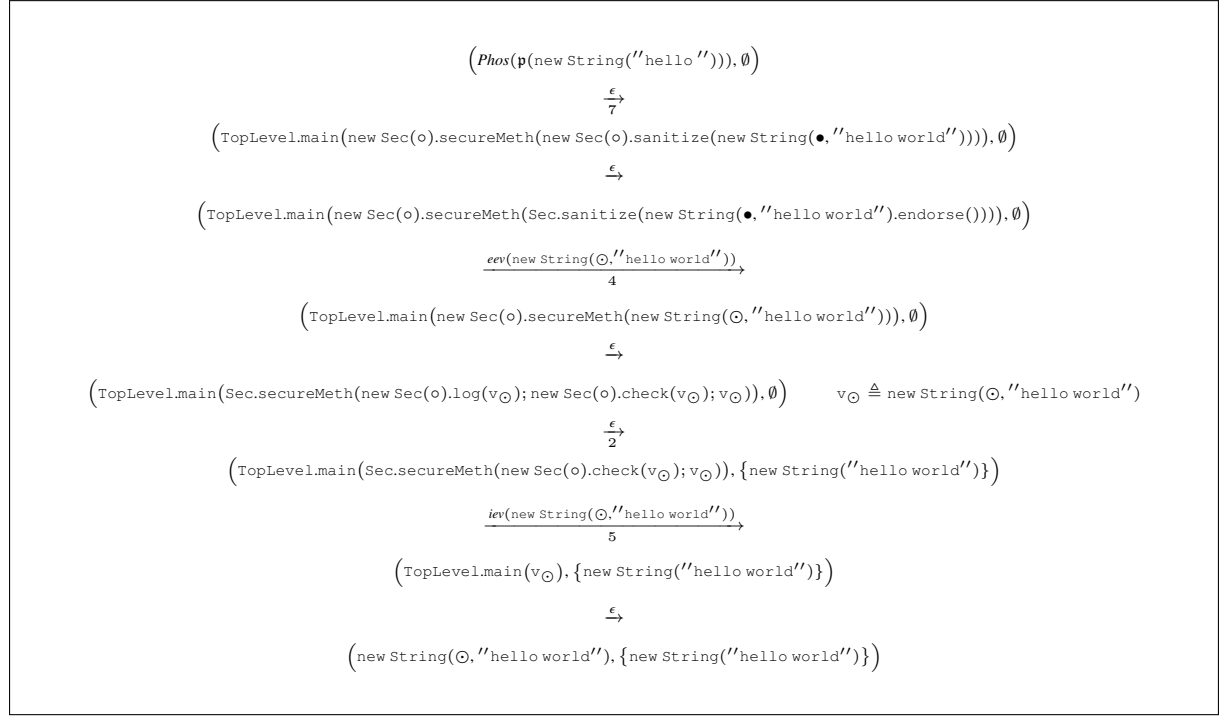
$$\left(\mathit{Phos}(\mathfrak{p}(\texttt{new String}(''\texttt{hello}''))), \emptyset\right)$$

$$\xrightarrow[7]{\epsilon}$$

$$\left(\texttt{TopLevel.main}(\texttt{new Sec}(\texttt{o}).\texttt{secureMeth}(\texttt{new Sec}(\texttt{o}).\texttt{sanitize}(\texttt{new String}(\bullet, ''\texttt{hello world}'')))), \emptyset\right)$$

$$\xrightarrow{\epsilon}$$

$$\left(\texttt{TopLevel.main}(\texttt{new Sec}(\texttt{o}).\texttt{secureMeth}(\texttt{Sec.sanitize}(\texttt{new String}(\bullet, ''\texttt{hello world}'').\texttt{endorse}()))), \emptyset\right)$$

$$\xrightarrow{\mathit{eev}(\texttt{new String}(\odot, ''\texttt{hello world}''))}$$

$$\xrightarrow{\quad\quad\quad 4 \quad\quad\quad}$$

$$\left(\texttt{TopLevel.main}(\texttt{new Sec}(\texttt{o}).\texttt{secureMeth}(\texttt{new String}(\odot, ''\texttt{hello world}''))), \emptyset\right)$$

$$\xrightarrow{\epsilon}$$

$$\left(\texttt{TopLevel.main}(\texttt{Sec.secureMeth}(\texttt{new Sec}(\texttt{o}).\texttt{log}(v_\odot); \texttt{new Sec}(\texttt{o}).\texttt{check}(v_\odot); v_\odot)), \emptyset\right) \quad\quad v_\odot \triangleq \texttt{new String}(\odot, ''\texttt{hello world}'')$$

$$\xrightarrow[2]{\epsilon}$$

$$\left(\texttt{TopLevel.main}(\texttt{Sec.secureMeth}(\texttt{new Sec}(\texttt{o}).\texttt{check}(v_\odot); v_\odot)), \{\texttt{new String}(''\texttt{hello world}'')\}\right)$$

$$\xrightarrow{\mathit{iev}(\texttt{new String}(\odot, ''\texttt{hello world}''))}$$

$$\xrightarrow{\quad\quad\quad 5 \quad\quad\quad}$$

$$\left(\texttt{TopLevel.main}(v_\odot), \{\texttt{new String}(''\texttt{hello world}'')\}\right)$$

$$\xrightarrow{\epsilon}$$

$$\left(\texttt{new String}(\odot, ''\texttt{hello world}''), \{\texttt{new String}(''\texttt{hello world}'')\}\right)$$

Fig. 12. Example 3.2: Target Trace.

and 4.2.3, and main proof details are presented in Section 4.2.4. Examples are provided to illustrate proof methodologies. *We note that throughout this Section we will ignore transition labels $\alpha$ in target language reduction since they are irrelevant to the properties of interest, and will use $\rightarrow$ exclusively to refer to the reduction relation in $FJ_{taint}$ defined in Section 4.*

Beginning with our main result for retrospective security, we note that our definition is general with respects to *SSOs* and *Sanitizers* defined at the top-level, which fix $LS_{taint}$ and $SP_{taint}$. Soundness and completeness as defined in Section 2.5 require definition of the notation $\tau \rightsquigarrow \mathbb{L}$, which for $FJ_{taint}$ means that $\mathbb{L}$ is the log in the last configuration of $\tau$. We define $toFOL(\mathbb{L}) = \{\text{MaybeBAD}(v) \mid v \in \mathbb{L}\}$, and thus $\lfloor \mathbb{L} \rfloor = C(toFOL(\mathbb{L}))$. Also as required, we define the relation $\cong$ below in Definition 4.7 establishing a semantic correspondence between FJ and $FJ_{taint}$ traces. Intuitively, the relation holds on source, target trace pairs if the taint shadow of configurations in the source trace match up with the structure of configurations in the target trace modulo security instrumentation.

**Theorem 4.1.** *For all $\mathfrak{p}(\mathbf{a})$, SSOs, and Sanitizers, let $\mathcal{R}(\mathfrak{p}(\mathbf{a}), LS_{taint}) = Phos(\mathfrak{p}(\mathbf{a}))$. Then $\mathcal{R}$ is semantics preserving, sound, and complete up to $SP_{taint}$.*

Since the safety property $SP_{taint}$ has been defined for FJ, operational correctness for prospective security means that any rewritten unsafe programs are blocked by instrumentation. We can formalize this property as follows, noting it is a consequence of semantics preservation under $\cong$. This is since bad shadows in source code correspond to values that fail security checks in the target.

$$overlay(\mathtt{x}, \mathtt{x}) = \mathtt{x} \qquad overlay(v, \delta) = v \qquad overlay(Op(\overline{\mathtt{e}}), Op(\overline{se})) = Op(\overline{overlay(\mathtt{e}, se)})$$

$$overlay(\mathtt{e.f}, se.\mathtt{f}) = overlay(\mathtt{e}, se).\mathtt{f} \qquad overlay(\mathtt{new}\ \mathtt{C}(\overline{\mathtt{e}}), \mathtt{shadow}\ \mathtt{C}(t, \overline{se})) = \mathtt{new}\ \mathtt{C}(t, \overline{overlay(\mathtt{e}, se)})$$

$$overlay(\mathtt{C.m}(\mathtt{e}), \mathtt{C.m}(se)) = \mathtt{C.m}(overlay(\mathtt{e}, se)) \qquad overlay(\mathtt{e.m}(\overline{\mathtt{e}'}), se.\mathtt{m}(\overline{se'})) = overlay(\mathtt{e}, se).\mathtt{m}(\overline{overlay(\mathtt{e}', se')})$$

$$overlay(\mathtt{if}\ \mathtt{e}_1\ \mathtt{then}\ \mathtt{e}_2\ \mathtt{else}\ \mathtt{e}_3, \mathtt{if}\ se_1\ \mathtt{then}\ se_2\ \mathtt{else}\ se_3) = \mathtt{if}\ overlay(\mathtt{e}_1, se_1)\ \mathtt{then}\ overlay(\mathtt{e}_2, se_2)\ \mathtt{else}\ overlay(\mathtt{e}_3, se_3)$$

Fig. 13. Definition of *overlay*.

**Definition 4.2.** *An FJ*$_{\mathrm{taint}}$ *program* $\mathtt{e}$ *causes a security failure iff*

$$(\mathtt{e}, \emptyset) \rightarrow^* (\mathtt{E}[\mathtt{v.check}(\mathtt{new}\ \mathtt{C}(\bullet, \overline{\mathtt{v}}))], \mathbb{L})$$

*for some* $\mathtt{E}$, $\mathtt{v}$, $\mathtt{new}\ \mathtt{C}(\bullet, \overline{\mathtt{v}})$, *and* $\mathbb{L}$.

Operational correctness of the prospective component of *Phos* can then be stated as follows:

**Theorem 4.2.** *The FJ program* $\mathfrak{p}(\mathbf{a})$ *is unsafe iff Phos*$(\mathfrak{p}(\mathbf{a}))$ *causes a security failure.*

*4.2.2. Preliminary Results for Semantic Correspondence*

Intuitively, we obtain a correspondence between source (FJ) and target (FJ$_{\mathrm{taint}}$) traces by "overlaying" shadows onto expressions in the former model, and by ignoring stuttering due to instrumentation steps in the latter. We begin with the relevant definitions for the source language, including the *overlay* function and a predicate for obtaining the last shadow in a trace.

**Definition 4.3.** *Given* $\tau$ *of length n, define*

$$\mathrm{LastShadow}(\tau) = se \iff \lfloor \tau \rfloor \otimes X \vdash \mathrm{LShadow}(se)$$

*where X contains the rules given in Figure 5, Figure 6 and the formula (1) as follows:*

$$\mathrm{Shadow}(n, se) \implies \mathrm{LShadow}(se). \tag{1}$$

*We define the "overlaying" of shadows on expressions in Figure 13.*

We observe that LastShadow is total function on nonempty traces, since shadow configurations are defined uniquely for every step of reduction in Figure 6.

For the target language, we define a relation $\hookrightarrow$ that is defined entirely in terms of $\rightarrow$ but elides instrumentation steps. This definition will eliminate stuttering and allow a much cleaner correspondence with source language reduction in proofs. In Definition 4.7 we will specify how to target traces with elided target traces.

**Definition 4.4.** *We define the operational specification relation* $\hookrightarrow$ *in terms of* $\rightarrow$, *in particular we define Phos*$(CT) \vdash \kappa \hookrightarrow \kappa'$ *if Phos*$(CT) \vdash \kappa \rightarrow \kappa'$[4] *for the* Sequence, Field, IfT, IfF, Return, *and* Var

---

[4]Recall from Section 3.1.4 that we write $CT \vdash \kappa \rightarrow \kappa'$ to be explicit about the class table $CT$ used for a reduction, and we write $CT \vdash_{\rightarrow} \tau$ to be explicit about the class table $CT$ and reduction relation $\rightarrow$ used in a trace $\tau$.

ElideInvoke

$$Phos(CT) \vdash (\texttt{new } \texttt{C}(\overline{v}).\texttt{m}(\overline{u}), \mathbb{L}) \rightarrow^* (\texttt{C.m}(\texttt{e}[\texttt{new } \texttt{C}(\overline{v})/\texttt{this}][\overline{u}/\overline{x}]), \mathbb{L}')$$

$$\frac{\mathit{mbody}_{CT}(\texttt{m}, \texttt{C}) = \overline{x}, \texttt{e}' \qquad \mathit{trim}(\texttt{e}) = \mu(\texttt{e}') \qquad \texttt{C.m} \notin \mathit{LibMeths}}{Phos(CT) \vdash (\texttt{new } \texttt{C}(\overline{v}).\texttt{m}(\overline{u}), \mathbb{L}) \hookrightarrow (\texttt{C.m}(\texttt{e}[\texttt{new } \texttt{C}(\overline{v})/\texttt{this}][\overline{u}/\overline{x}]), \mathbb{L}')}$$

ElideProp

$$\frac{Phos(CT) \vdash (\texttt{new } \texttt{C}(\overline{v}).\texttt{m}(\overline{u}), \mathbb{L}) \rightarrow^* (\texttt{C.m}(\texttt{new } \texttt{D}(t, \overline{e})), \mathbb{L}') \qquad \mathit{mbody}_{CT}(\texttt{m}, \texttt{C}) = \overline{x}, \texttt{new } \texttt{D}(\overline{e}') \qquad \overline{e} = \mu(\overline{e}') \qquad \texttt{C.m} \in \mathit{LibMeths}}{Phos(CT) \vdash (\texttt{new } \texttt{C}(\overline{v}).\texttt{m}(\overline{u}), \mathbb{L}) \hookrightarrow (\texttt{C.m}(\texttt{new } \texttt{D}(t, \overline{e})), \mathbb{L}')}$$

ElideEndorse

$$\frac{Phos(CT) \vdash (\texttt{C.m}(\texttt{v.endorse}()), \mathbb{L}) \rightarrow^* (\texttt{C.m}(\texttt{u}), \mathbb{L})}{Phos(CT) \vdash (\texttt{C.m}(\texttt{v}), \mathbb{L}) \hookrightarrow (\texttt{u}, \mathbb{L})}$$

ElideContext

$$\frac{Phos(CT) \vdash (\texttt{e}, \mathbb{L}) \hookrightarrow (\texttt{e}', \mathbb{L}')}{Phos(CT) \vdash (\texttt{E}[\texttt{e}], \mathbb{L}) \hookrightarrow (\texttt{E}[\texttt{e}'], \mathbb{L}')}$$

Fig. 14. Operational Specification Relation $\hookrightarrow$ to Elide Instrumentation Steps.

$$(Phos(\mathfrak{p}(\texttt{new String}(''\texttt{hello}'')))), \emptyset)$$

$$\hookrightarrow^5$$

$$(\texttt{TopLevel.main}(\texttt{new Sec}(\texttt{o}).\texttt{secureMeth}(\texttt{new Sec}(\texttt{o}).\texttt{sanitize}(\texttt{new String}(\bullet, ''\texttt{hello world}'')))), \emptyset)$$

$$\hookrightarrow^2$$

$$(\texttt{TopLevel.main}(\texttt{new Sec}(\texttt{o}).\texttt{secureMeth}(\texttt{new String}(\odot, ''\texttt{hello world}''))), \emptyset)$$

$$\hookrightarrow^2$$

$$(\texttt{TopLevel.main}(\texttt{new String}(\odot, ''\texttt{hello world}'')), \{\texttt{new String}(''\texttt{hello world}'')\})$$

$$\hookrightarrow$$

$$(\texttt{new String}(\odot, ''\texttt{hello world}''), \{\texttt{new String}(''\texttt{hello world}'')\})$$

Fig. 15. Example 3.2: Corresponding $\hookrightarrow$ trace.

*cases, but redefine* Invoke *as* ElideInvoke *and* ElideProp *allowing security instrumentation steps to be skipped during reduction, and add an* ElideEndorse *rule that elides explicit steps for endorsing a value in Figure 14. Note in the definition of* ElideInvoke *that the class table CT is assumed to be the one prior to rewriting, while $\rightarrow$ reductions assume the class table Phos(CT). Also in this definition trim is a function on expressions* $\texttt{e}$ *that replaces all subexpressions of* $\texttt{e}$ *of the form* $\texttt{e}'.\texttt{endorse}()$ *with* $\texttt{e}'$. *As for $\rightarrow$ we will write* $\kappa \hookrightarrow \kappa'$ *when CT is clear from context.*

**Example 4.1.** *An example $\hookrightarrow$ trace $\tau$ is shown in Figure 15, which extends the running Example 3.2– in particular $\mathfrak{p}$ is as defined there. Note how $\sigma$ reduces in lockstep with the source reduction in Figure 8.*

We break down the proof of semantic correspondence into two main lemmas for both the soundness and completeness directions. Given the above definitions they follow in a straightforward manner by case analysis on the form of reduction. In the following we assume a fixed set of *SSOs* and *Sanitizers*, though the proofs are abstract with regard to their contents beyond the sanity conditions (closure under inheritance) described in Section 3.2.3.

**Lemma 4.1.** *Assume given a source trace $CT \vdash_{\rightarrow} \sigma$ and elided target trace $Phos(CT) \vdash_{\hookrightarrow} \tau$ that end with configurations $\mathsf{e}_s$ and $(\mathsf{e}_t, \mathbb{L})$ respectively where by assumption the following hold:*

$$\mathrm{LastShadow}(\sigma) = se \qquad overlay(se, \mathsf{e}_s) = \mathsf{e}_t \qquad \mathsf{e}_s \rightarrow \mathsf{e}_s' \qquad \mathrm{LastShadow}(\sigma \mathsf{e}_s') = se'$$

$$\sigma \mathsf{e}_s' \in \mathrm{SP}_{\mathrm{taint}}$$

*Then there exists a target configuration $(\mathsf{e}_t', \mathbb{L}')$ where the following hold:*

$$(\mathsf{e}_t, \mathbb{L}) \hookrightarrow (\mathsf{e}_t', \mathbb{L}') \qquad\qquad overlay(se', \mathsf{e}_s') = \mathsf{e}_t'$$

**Proof.** By case analysis on $\mathsf{e}_s \rightarrow \mathsf{e}_s'$. Given the close correspondence of $\rightarrow$ and $\hookrightarrow$ reduction the proof proceeds trivially in all cases except for (i) the subcase of case Invoke where an sso is invoked with an untainted or maybe tainted argument or (ii) the subcase in which a library method is invoked, or (iii) in the subcase of case Return where a sanitizer has completed execution.

For (i) we have by assumption that the shadow of the sso's argument will reflect its taint level, which will also be the taint level of its corresponding target object since $overlay(se, \mathsf{e}_s) = \mathsf{e}_t$ by assumption. Thus, the desired target reduction exists since in particular the elided `check` will succeed.

For (ii) we likewise have by assumption that the shadow of the library method's arguments will reflect its taint level, which will also be the taint level of its corresponding target objects since $overlay(se, \mathsf{e}_s) = \mathsf{e}_t$ by assumption. The result follows by definition of Prop which will propagate taint in the same way as the instrumentation in the library method by assumption.

For (iii) we have by assumption that the returned value will be endorsed, as reflected in its shadow. Furthermore the desired target reduction exists by definition of $\rightarrow$, which is a return from a sanitizer, hence by definition of Endorse the returned object is properly endorsed in the target, due to the invocation of `endorse` in the elided steps. □

**Lemma 4.2.** *Assume given a source trace $CT \vdash_{\rightarrow} \sigma$ and elided target trace $Phos(CT) \vdash_{\hookrightarrow} \tau$ that end with configurations $\mathsf{e}_s$ and $(\mathsf{e}_t, \mathbb{L})$ respectively where by assumption the following hold:*

$$\mathrm{LastShadow}(\sigma) = se \qquad overlay(se, \mathsf{e}_s) = \mathsf{e}_t \qquad \sigma \in \mathrm{SP}_{\mathrm{taint}} \qquad (\mathsf{e}_t, \mathbb{L}) \hookrightarrow (\mathsf{e}_t', \mathbb{L}')$$

*Then there exists source and shadow expressions $\mathsf{e}_s'$ and $se'$ where the following hold:*

$$\mathsf{e}_s \rightarrow \mathsf{e}_s' \qquad \mathrm{LastShadow}(\sigma \mathsf{e}_s') = se' \qquad overlay(se', \mathsf{e}_s') = \mathsf{e}_t' \qquad \sigma \mathsf{e}_s' \in \mathrm{SP}_{\mathrm{taint}}$$

**Proof.** By case analysis on $(\mathsf{e}_t, \mathbb{L}) \hookrightarrow (\mathsf{e}_t', \mathbb{L}')$. Given the close correspondence of $\rightarrow$ and $\hookrightarrow$ reduction the proof proceeds trivially in all cases except for cases (i) ElideEndorse, (ii) ElideInvoke, and (iii) ElideProp.

For (i) we have by assumption that the result is endorsed, and by definition of Shadow we have also that the returned value of a sanitizer is endorsed as reflected in its shadow, hence the result follows.

For (ii) we have by assumption that the argument of the sso is either untainted or maybe tainted, since otherwise the assumed $\hookrightarrow$ relation would not hold. The source configuration $\mathsf{e}_s$ will thus not be BAD since by assumption $overlay(se, \mathsf{e}_s) = \mathsf{e}_t$, so the desired result follows by definition.

For (iii), since by assumption $overlay(se, \mathsf{e}_s) = \mathsf{e}_t$ the target and source configurations $\mathsf{e}_t$ and $\mathsf{e}_s$ will both be invocations of the same library method on arguments with corresponding taint labels. The result

follows by assumed correspondence of Prop and the instrumentation of taint propagation in the library method. □

The above Lemmas will allow an induction step for our proof of the main Theorem 4.1. The following Lemma establishes the basis of the induction, and follows trivially by definition.

**Lemma 4.3.** *For all FJ top-level programs* $\mathfrak{p}(\mathbf{a})$ *with* $\mathbf{a}$ *an object in class* C, *letting:*

$$se = \texttt{shadow TopLevel(o).main(shadow C}(\bullet, \delta))$$

*it is the case that:*

$$overlay(se, \mathfrak{p}(\mathbf{a})) = Phos(\mathfrak{p}(\mathbf{a}))$$

Given the last few Lemmas, we can show any top-level program source trace has an elided target trace of the same length, with corresponding configurations proceeding in lockstep. We define the trace correspondence relation $\simeq$ appropriately and show that it holds for any program and attack.

**Definition 4.5.** *For all* $CT \vdash_{\rightarrow} \sigma$ *and* $Phos(CT) \vdash_{\hookrightarrow} \tau$, *define* $\sigma \simeq \tau$ *as the least relation satisfying the following inductive rules:*

$$\varnothing \simeq \varnothing \qquad \qquad \frac{overlay(\mathsf{e}_s, \mathrm{LastShadow}(\sigma\mathsf{e}_s)) = \mathsf{e}_t}{\sigma\mathsf{e}_s \simeq \tau(\mathsf{e}_t, \mathbb{L})}$$

**Lemma 4.4.** *For all* $\mathfrak{p}$ *and* $\mathbf{a}$ *there exists* $\sigma \in \mathrm{SP}_{\mathrm{taint}}$ *with* $CT \vdash_{\rightarrow} \mathfrak{p}(\mathbf{a})\sigma$ *iff there exists* $\tau$ *with* $Phos(CT) \vdash_{\rightarrow} (Phos(\mathfrak{p}(\mathbf{a}), \varnothing)\tau$ *and* $\mathfrak{p}(\mathbf{a})\sigma \simeq (Phos(\mathfrak{p}(\mathbf{a}), \varnothing)\tau$

**Proof.** Straightforward by induction on the length of traces and applications of Lemma 4.3 for the base (length 1) case and Lemmas 4.1 and 4.2 in the inductive cases. □

### 4.2.3. Preliminary Results for Logging Correctness

To establish correctness of logging in the image of *Phos* some steps must be taken to show that adding tainted or maybe tainted values to the log obtains the same information as resolution of the MaybeBAD predicate. Our main argument (Lemma 4.10) focuses on steps that change the log/logging specification.

Since instrumented logs and logging specifications have different representations, we must show these represent the same information. As in [8], our proof strategy is based on an appeal to least Herbrand models $\mathfrak{H}$ of the logging specifications and logs (least Herbrand models are known to exist for safe Horn clause logics with compound terms [49]). In essence, we demonstrate that audit logs generated by FJ$_{\mathrm{taint}}$ programs *are* the least Herbrand model of the logging specification for the source program, hence contain the same information. We start by demonstrating the following implication that allows us to consider logging specifications in terms of least Herbrand models $\mathfrak{H}$ of sets of formulas.

**Lemma 4.5.** *Given X as defined for* $LS_{\mathrm{taint}}$, *for all* $\mathbb{L}$ *and* $\tau\kappa$, *if:*

$$toFOL(\mathbb{L}) = \mathfrak{H}(X \cup toFOL(\tau)) \cap L_{\{\mathrm{MaybeBAD}\}}$$

*then* $\lfloor \mathbb{L} \rfloor = LS_{\mathrm{taint}}(\tau\kappa)$.

**Proof.** By assumption and closure properties of least Herbrand models [49] we have:

$$C(toFOL(\mathbb{L})) = C(C(X \cup toFOL(\tau)) \cap L_{\{\text{MaybeBAD}\}})$$

so the result follows by Definition. □

Next, we establish a property that can be applied for our inductive argument in Lemma 4.14, which shows that adding a value $v$ to a log is the same as adding a new resolution MaybeBAD($v$) to the logging specification.

**Lemma 4.6.** *Given $X$ as defined for $LS_{\text{taint}}$ and $\mathbb{L}$ and $\tau$ such that:*

$$toFOL(\mathbb{L}) = \mathfrak{H}(X \cup toFOL(\tau)) \cap L_{\{\text{MaybeBAD}\}}$$

*if* C.m $\in SSOs$ *and* LastShadow($\tau$) = $SE[\text{shadow C}(t', \overline{se}).\text{m}(\text{shadow D}(t, \overline{\delta}))]$ *for some* D, $\overline{se}$, $t'$, *and* $t \leqslant \odot$ *it is the case that:*

$$toFOL(\mathbb{L} \cup v) = \mathfrak{H}(X \cup toFOL(\tau E[\text{new C}(\overline{v}).\text{m}(v)])) \cap L_{\{\text{MaybeBAD}\}}$$

**Proof.** The result follows by observing that adding $v$ to the $\mathbb{L}$ is the same as adding MaybeBAD($v$) to $toFOL(\mathbb{L})$ by definition, while adding $E[\text{new C}(\overline{v}).\text{m}(v)]$ to the trace $\tau$ adds a new resolution MaybeBAD($v$) to $X \cup toFOL(\tau)$. Thus, the change in the least Herbrand model with respect to $L_{\{\text{MaybeBAD}\}}$ will therefore also just be the addition of MaybeBAD($v$). □

The following result establishes the base case for the induction and holds immediately by definition.

**Lemma 4.7.** *Given $X$ as defined for $LS_{\text{taint}}$, for all $\mathfrak{p}$ and $\mathbf{a}$ we have:*

$$toFOL(\emptyset) = \mathfrak{H}(X \cup toFOL(\varnothing)) \cap L_{\{\text{MaybeBAD}\}} = C(\varnothing)$$

The next Lemma shows that if a source trace encounters an sso invocation with a maybe tainted argument, then that argument will be added to the log in the corresponding target trace.

**Lemma 4.8.** *Assume given a source trace $CT \vdash_{\rightarrow} \sigma e_s$ and elided target trace $Phos(CT) \vdash_{\hookrightarrow} \tau$ that ends with configuration $(e_t, \mathbb{L})$ where the following hold:*

$$\sigma e_s \simeq \tau \qquad \text{LastShadow}(\sigma) = SE[\text{shadow C}(t', \overline{se}).\text{m}(\text{shadow D}(t, \overline{\delta}))] \qquad \text{C.m} \in SSOs \qquad t \leqslant \odot$$

*Then $v \in \mathbb{L}$.*

**Proof.** By Definition of $\simeq$ and *overlay*, it must be the case that the penultimate configuration in $\tau$ is a call to C.m with the same argument explicitly labeled with $t$. □

The following Lemma shows that if a value is added to the log in a target trace, then its corresponding source trace must have a call to an sso with the same value that is maybe tainted.

$$Phos(\mathfrak{p}(\mathbf{a})) \rhd Phos(\mathfrak{p}(\mathbf{a})) \qquad \frac{\tau \rhd \tau'}{\tau \kappa_2 \rhd \tau' trim(\kappa_2)} \qquad \frac{\tau \rhd \tau' \qquad e' = \text{u.log(v)} \vee e' = e''; \text{u.log(v)}}{\tau(\text{E}[e'; e]), \mathbb{L}) \rhd \tau'}$$

$$\frac{\tau \rhd \tau'}{\tau(\text{E}[v; e]), \mathbb{L}) \rhd \tau'(\text{E}[e], \mathbb{L})} \qquad \frac{\tau \rhd \tau' \kappa_1 \qquad \neg(\kappa_2 = \text{E}[\text{C.m}(e)] \wedge \text{C.m} \in SSOs) \qquad \neg(\kappa_1 \hookrightarrow trim(\kappa_2))}{\tau \kappa_2 \rhd \tau' \kappa_1}$$

Fig. 16. Definition of $\rhd$ Relating Target and Elided Target Traces.

**Lemma 4.9.** *Assume given a source trace* $CT \vdash_{\to} \sigma e_s$ *and elided target trace* $Phos(CT) \vdash_{\hookrightarrow} \tau$ *that ends with the trace* $(e_{t_1}, \mathbb{L}_1)(e_{t_2}, \mathbb{L}_2)$ *with* $\sigma e_s \simeq \tau$, $v \notin \mathbb{L}_1$, *and* $v \in \mathbb{L}_2$. *Then the following hold:*

$$\text{LastShadow}(\sigma) = SE[\text{shadow C}(t', \overline{se}).\text{m}(\text{shadow D}(t, \overline{\delta}))] \qquad \text{C.m} \in SSOs \qquad t \leqslant \odot$$

**Proof.** If $v$ is added to the log in the last step of $\tau$, then by definition of $\hookrightarrow$ its penultimate configuration is a call to an sso with $v'$ as an argument, where $v'$ is the same as $v$ except with explicit taint label $t \leqslant \odot$. Thus the result follows by definition of $\simeq$. $\qquad \square$

Now we can demonstrate a central Lemma, establishing that logging is correct according to $\hookrightarrow$.

**Lemma 4.10.** *Let X be defined as for* $LS_{\text{taint}}$. *If* $\sigma \simeq (e, \varnothing)\tau$ *where* $(e, \varnothing)\tau \rightsquigarrow \mathbb{L}$ *then:*

$$toFOL(\mathbb{L}) = \mathfrak{H}(X \cup toFOL(\sigma)) \cap L_{\{\text{MaybeBAD}\}}$$

**Proof.** By induction on the length $n$ of $\tau$. In case $n = 1$, the result follows by the Definition of $\simeq$ and Lemma 4.7. In case $n > 1$, the result follows trivially by the induction hypothesis except in case the source or target trace adds information to the logging specification of log, respectively.

Suppose on the one hand that information in $LS_{\text{taint}}(\sigma)$ is changed at the end of $\sigma$– by definition of $LS_{\text{taint}}$ and MaybeBAD this means that the penultimate configuration in $\sigma$ is an invocation of an sso on an argument $v$ with shadow taint $t \leqslant \odot$. But then by Lemma 4.8 $v \in \mathbb{L}$. Suppose on the other hand that a value $v$ is added to the log in the last step of $(e, \varnothing)\tau$. Then by Lemma 4.9 the last step in $\sigma$ is an invocation of an sso on an argument $v$ with shadow taint $t \leqslant \odot$. The desired result then follows by Lemma 4.6. $\qquad \square$

### 4.2.4. Definition of $\cong$ and Proofs of Main Results

The above preliminary results do most of the heavy lifting, but so far we have only related source traces with elided target traces, and we must relate them to plain target traces with $\cong$. To accomplish this we define the relation $\rhd$ that associates target traces with their corresponding elided traces. An important subtlety of the definition of $\tau \rhd \tau'$ is that when $\tau$ is within an sso, $\tau'$ must end *before* information is added to the log if the relevant call to $\log$ has not yet been evaluated in $\tau$, and *after* otherwise. To accommodate this we take advantage of the fact that sequencing syntax using ; only occurs in instrumentation steps in the image of *Phos*, specifically within an sso, and thus use that symbol as a marker. We clarify that ; is left-associative, so in particular a target redex $\text{E}[v; e]$ is a computation currently within an sso where a check has just completed.

**Definition 4.6.** *For all $\tau$ and $\tau'$ where $Phos(CT) \vdash_{\rightarrow} \tau$ and $Phos(CT) \vdash_{\hookrightarrow} \tau'$, we say $\tau \rhd \tau'$ iff the relation can be derived using the rules defined in Figure 16, where abusing notation we define $trim(e, \mathbb{L}) = (trim(e), \mathbb{L})$.*

It is easy to show that $\rhd$ is a surjective function from $\rightarrow$ traces to $\hookrightarrow$ traces by induction on the length of traces and the definition of $\hookrightarrow$, which is entirely in terms of $\rightarrow$. Also, the definition of $\rhd$ guarantees that related traces will always generate the same log.

**Lemma 4.11.** *Given source CT each of the following conditions hold:*

(1) *For all $\tau$ where $Phos(CT) \vdash_{\rightarrow} \tau$ there exists unique $\sigma$ where $Phos(CT) \vdash_{\hookrightarrow} \sigma$ and $\tau \rhd \sigma$.*
(2) *For all $\sigma$ where $Phos(CT) \vdash_{\hookrightarrow} \sigma$ there exists $\tau$ where $Phos(CT) \vdash_{\rightarrow} \tau$ and $\tau \rhd \sigma$.*
(3) *For all $\tau$ where $Phos(CT) \vdash_{\rightarrow} \tau$ and $\sigma$ where $Phos(CT) \vdash_{\hookrightarrow} \sigma$, if $\tau \rhd \sigma$ then $\tau \rightsquigarrow \mathbb{L}$ iff $\sigma \rightsquigarrow \mathbb{L}$.*

Now, given the above we can define $\cong$ as follows. The Definition is in terms of $\rhd$ and $\simeq$ allowing us to take advantage of preliminary results.

**Definition 4.7** (FJ and FJ$_{taint}$ Correspondence Relation). *Given source language execution trace $\sigma$ and target language execution trace $\tau$ define:*

$$\sigma \cong \tau \iff \exists \tau'. \tau \rhd \tau' \wedge \sigma \simeq \tau'$$

**Example 4.2.** *Letting $\sigma$ be the full trace shown in Figure 8, $\tau$ be the full trace shown in Figure 12, and $\tau'$ be the full trace shown in Figure 15, we have $\tau \rhd \tau'$, $\sigma \simeq \tau'$, and $\sigma \cong \tau$.*

Now we can prove our main results. The following three Lemmas establish the necessary conditions for semantic preservation and logging correctness with respect to $\cong$, and straightforward proofs of the main Theorems follow.

**Lemma 4.12** (Simulation Soundness of *Phos*). *If $\mathfrak{p}(\mathbf{a}) \Downarrow \sigma$ and $\sigma \in \mathrm{SP}_{taint}$ then there exists $\tau$ where $Phos(\mathfrak{p}(\mathbf{a})) \Downarrow \tau$ with $\sigma \cong \tau$.*

**Proof.** By Lemma 4.4 there exists $\tau'$ such that $\sigma \simeq \tau'$, and by Lemma 4.11 there exists $\tau$ such that $\tau \rhd \tau'$. The result follows by definition of $\cong$. $\square$

**Lemma 4.13** (Simulation Completeness of *Phos*). *If $Phos(\mathfrak{p}(\mathbf{a})) \Downarrow \tau$ then there exists $\sigma \in \mathrm{SP}_{taint}$ where $\mathfrak{p}(\mathbf{a}) \Downarrow \sigma$ and $\sigma \cong \tau$.*

**Proof.** By Lemma 4.11 there exists $\tau'$ such that $\tau \rhd \tau'$, and by Lemma 4.4 there exists $\sigma$ such that $\sigma \simeq \tau'$. The result follows by definition of $\cong$. $\square$

**Lemma 4.14** (Logging Correctness of *Phos*). *For all $\mathfrak{p}$ and $\mathbf{a}$, if $\mathfrak{p}(\mathbf{a}) \Downarrow \sigma$ and $Phos(\mathfrak{p}(\mathbf{a})) \Downarrow \tau$ with $\sigma \cong \tau$ and $\tau \rightsquigarrow \mathbb{L}$ then $\lfloor \mathbb{L} \rfloor = LS_{taint}(\sigma)$*

**Proof.** By definition of $\cong$ there exists $\tau'$ where $\tau \rhd \tau'$ and $\sigma \simeq \tau'$. But by Lemma 4.11 $\tau' \rightsquigarrow \mathbb{L}$, and therefore the result follows by Lemma 4.10 and Lemma 4.5. $\square$

*Proofs of main results.*    Now we can easily demonstrate the main results on the basis of the preceding lemmas.

**Proof of Theorem 4.1.** Immediate by Lemmas 4.12, 4.13, and 4.14.                                  □

**Proof of Theorem 4.2.** The result follows directly by Lemmas 4.12 and 4.13 and the definition of $\cong$, which together guarantee that a BAD configuration is encountered upon invocation of an sso in the source iff a target sso is invoked with a tainted argument.                                  □

## 5. The Security Property of *Phos*

The semantics of information flow has been well studied and is typically characterized via noninterference properties, but surprisingly little work has been done to develop similar properties for taint analysis. In recent years it has been shown that direct flow of data confidentiality is not comparable with noninterference [10], i.e., there are both noninterfering programs with direct leakage of secret data to public domain, and programs without such direct leakages, but interfering. For instance consider the following two statements in a core imperative language, in which $s$ and $p$ are respectively secret and public variables:

$$\textbf{if } s = 0 \textbf{ then } p := s \textbf{ else } p := 0 \qquad\qquad \textbf{if } s = 0 \textbf{ then } p := 1 \textbf{ else } p := 2$$

The first statement is noninterfering, but direct flow of information from $s$ to $p$ exists, whereas the second statement is interfering due to the indirect flow from $s$ to $p$, but there are not any direct flows from $s$ to $p$.

Formal definitions of taint analysis implementations do exist, but they are usually operational in nature. For example, in Section 4.2, we have established an operational correctness result for the prospective enforcement of direct integrity flow. In this section, we propose a hyperproperty to characterize the security property enforced by integrity taint analysis techniques. This hyperproperty is defined in a general, language-agnostic way, though in this Section we also show that the instrumentation of $FJ_{taint}$ programs by *Phos* enjoys this property as a correctness condition. We illustrate key points in Section 5.3.

### 5.1. Direct Integrity Flow Semantics: Explicit Integrity

We define *explicit integrity* as a semantic hyperproperty that builds on (dualizes) the notions of explicit secrecy [10] and attacker power [34]. Similar to explicit secrecy, explicit integrity is language-agnostic. In later sections, we discuss instantiation of this model for $FJ_{taint}$.

Intuitively, a program enjoys explicit secrecy if execution of its state transformation components does not affect the knowledge of a low confidentiality user. By formally specifying state transformation components, control flow operations (such as conditional expressions) can be omitted to only consider direct aka explicit program flows. *Knowledge* [30] is defined as the set of initial states configurable by a low confidentiality user that generate a particular sequence of observables– the smaller the set, the greater the knowledge. *Explicit knowledge* [10] restricts this concept to direct program data flow. In this section, we demonstrate how explicit knowledge can be "dualized" for direct integrity flow analysis and applied

as a semantic framework for dynamic integrity taint analysis tools, particularly in functional languages with hierarchical data structures (FJ$_{taint}$).

*Attacker power* [34] is introduced as a counterpart to attacker knowledge in the context of integrity, as the set of low integrity inputs that generate the same sequence of high integrity events. Each high integrity event could be a simple assignment to a predefined high integrity variable, a method that manipulates trusted data (secure sinks), etc. according to the language model. The more refined the attacker power is, the more powerful the low integrity attacker becomes, as she becomes more capable to distinguish between the effects of different attacks on high integrity data.

We define *explicit attacker power* as the attacker power constrained on direct integrity flows. Then, *explicit integrity* is defined as the property of preserving explicit attacker power during program execution. In order to limit flows to direct ones, we have followed the techniques introduced in [10] to define *state transformers*. State transformers extract direct flows semantically by specifying the ways in which program state is modified in each step of execution, along with direct-flow events that are generated.

### 5.1.1. Model Specification

We formulate our explicit integrity semantics following [10]. We first define the interface for our framework. Let $\mathbf{K}$ be the set of program configurations for a given object language where $\kappa$ ranges over configurations. Configurations consist of control and state segments. Control refers to code and state refers to data. Let $\mathbf{C}$ be the set of controls with $\mathbf{c}$ ranging over the elements of $\mathbf{C}$. Moreover, let $\mathbf{S}$ denote the set of states and $\mathbf{s}$ represent a given state. We also define a set of high integrity events, $\mathbf{E}$. A high integrity event $\mathbf{e}$ may refer to different computations in different language models and settings. For example, it could be as simple as assigning a low integrity data to a high integrity variable, or invoking a method with a low integrity data as its parameter to store that parameter in a database. We let $\alpha$ range over elements of $\mathbf{E}^*$. We assume the existence of the evaluation relation $\to \subseteq \mathbf{K} \times \mathbf{E}^* \times \mathbf{K}$ where $(\kappa, \alpha, \kappa') \in \to$ is denoted as $\kappa \xrightarrow{\alpha} \kappa'$. We use $\kappa \to \kappa'$ if $\alpha$ is empty ($\epsilon$) or could be elided in the discussion. Notation $\to^*$ is used for reflexive and transitive closure of $\to$.

Each configuration is considered to include two segments: control (code) and state (data). These segments are not necessarily disjoint and could overlap in some language models. In this regard, let mappings $state : \mathbf{K} \to \mathbf{S}$ and $com : \mathbf{K} \to \mathbf{C}$ extract the state and control segments of configurations, and $\langle \cdot, \cdot \rangle : \mathbf{C} \times \mathbf{S} \to \mathbf{K}$ construct a configuration from its control and state segments. These mappings need to satisfy the following property, for any $\kappa$:

$$\langle com(\kappa), state(\kappa) \rangle = \kappa.$$

We assume the existence of an entry point $[\cdot]$ in the controls denoted by $\mathbf{c}[\cdot]$ by which the attacker can inject low integrity input. The attacker input is denoted by $\mathbf{a}$. Then $\mathbf{c}[\mathbf{a}]$ represents a control in which the attacker has injected input $\mathbf{a}$. Note that an attack $\mathbf{a}$ is a data piece itself, i.e., $\mathbf{a}$ is a value.

We define extracted state transformers as follows. A consideration of state transformation, rather than complete program execution, allows us to focus only on direct program flow, rather than indirect control flow e.g. via conditionals. State transformers play the same role that explicit flow statements do in Weak Secrecy [35]. We note that this definition is a slight refinement of the analogous definition in [10]– in their work, a command is assumed to be compatible with all states, whereas we require compatibility of commands and states. This refinement is necessary due to structured expressions in HLLs such as Java, vs. lower level languages. However, we add a completeness condition expressed in Definition 5.3 that ensures we can compare all trust equivalent states via state transformation functions.

**Definition 5.1.** *Let $\kappa \to \kappa'$ and $com(\kappa) = \mathbf{c}$ for some $\mathbf{c}$. $f : \mathbf{S} \to \mathbf{S} \times \mathbf{E}^*$ is the function where $f(\mathbf{s}) = (state(\kappa''), \alpha)$ for all $\mathbf{s}$ such that $\langle \mathbf{c}, \mathbf{s} \rangle$ is defined and for the unique $\kappa''$ and $\alpha$ such that $\langle \mathbf{c}, \mathbf{s} \rangle \xrightarrow{\alpha} \kappa''$. We write $\kappa \to_f \kappa'$ to associate the state transformer $f$ with the reduction $\kappa \to \kappa'$. This definition is then extended to multiple evaluation steps by composing state transformers at each step. Let $f(\mathbf{s}) = (\mathbf{s}', \alpha)$ and $g(\mathbf{s}') = (\mathbf{s}'', \alpha')$. Then, $(g * f)(\mathbf{s}) = (\mathbf{s}'', \alpha\alpha')$.*

We now define the power an attacker obtains by observing high integrity events. We capture this by defining a set of high integrity equivalent states that generate the same sequence of high integrity events. We posit the binary relation $=_\circ$ on $\mathbf{S}$ to denote high integrity equivalent (or trust equivalent) states. The general sense of this relation is that $\mathbf{s} =_\circ \mathbf{s}'$ if $\mathbf{s}$ and $\mathbf{s}'$ agree on high integrity data. The instantiation of the relation depends on the language model in which the states are defined. For a state $\mathbf{s}$ and some state transformer $f$, the state $\mathbf{s}'$ is considered as an element of the explicit attacker power if $\mathbf{s} =_\circ \mathbf{s}'$ and $\mathbf{s}'$ agrees with $\mathbf{s}$ on the generated high integrity events.

**Definition 5.2.** *We define* explicit attacker power *with respect to state $\mathbf{s}$ and state transformer $f$ as follows, where projection on the ith element of a tuple is denoted by $\pi_i$.*

$$p_e(\mathbf{s}, f) = \{\mathbf{s}' \mid \mathbf{s} =_\circ \mathbf{s}', \pi_2(f(\mathbf{s})) = \pi_2(f(\mathbf{s}'))\}.$$

All state transformers must be complete in the following sense for this definition to be coherent:

**Definition 5.3.** *A state transformer $f$ is* complete *iff for all $\mathbf{s}_1$, $\mathbf{s}_2$ where $\mathbf{s}_1 =_\circ \mathbf{s}_2$ we have $f(\mathbf{s}_1)$ is defined iff $f(\mathbf{s}_2)$ is defined.*

A control then satisfies explicit integrity for some state iff no state can be excluded from observing the high integrity events generated by the extracted state transformer.

**Definition 5.4.** *A control $\mathbf{c}$ satisfies* explicit integrity *for state $\mathbf{s}$, iff $\langle \mathbf{c}, \mathbf{s} \rangle \to_f^* \kappa'$ implies that for any $\mathbf{s}'$ and $\mathbf{s}''$, if $\mathbf{s}' =_\circ \mathbf{s}''$ then we have $\mathbf{s}'' \in p_e(\mathbf{s}', f)$. A control $\mathbf{c}$ satisfies* explicit integrity *iff for any $\mathbf{s}$, $\mathbf{c}$ satisfies explicit integrity for $\mathbf{s}$.*

We can now consider explicit integrity in the presence of endorsement in the style of gradual release [30]. We assume that there exists a set of integrity events $\mathbf{E_{en}} \subseteq \mathbf{E}$ that are generated when endorsements occur. Explicit attacker power is only allowed to change for such events.

**Definition 5.5.** *A control $\mathbf{c}$ satisfies* explicit integrity modulo endorsement *for state $\mathbf{s}$ iff $\langle \mathbf{c}, \mathbf{s} \rangle \to_f^* \kappa' \xrightarrow{\alpha}_g^* \kappa''$ and $\alpha \notin \mathbf{E_{en}}^*$ imply that $p_e(\mathbf{s}, f) = p_e(\mathbf{s}, g * f)$.*

*5.2. An Instantiation with FJ*$_{taint}$

In this section, we instantiate explicit integrity for FJ$_{taint}$. *Because audit logging and retrospective features are irrelevant to the technical development in this Section, we omit them and elide FJ$_{taint}$ configurations to just expressions $\mathbf{e}$, and take* log *to be the identity function.*

First, we define the required interface specified in Section 5.1, beginning with the definition of extracted state transformers for all features. These are extracted from the definition of $\to$— notably, the extracted state transformers for conditional expressions inline conditional branching, disregarding the

$$\frac{fields_{Phos(CT)}(\text{C}) = \overline{\text{C}}\ \overline{\text{f}} \qquad \text{f}_i \in \overline{\text{f}}}{select_{\text{f}_i}(\text{E}[\text{new C}(\overline{v}).\text{f}_i]) = (\text{E}[v_i], \epsilon)} \qquad \frac{mbody_{Phos(CT)}(\text{m}, \text{C}) = \overline{\text{x}}, \text{e} \qquad \text{C.m} \notin LibMeths}{call_{\text{C.m}}(\text{E}[\text{new C}(\overline{v}).\text{m}(\overline{\text{u}})]) = (\text{E}[\text{C.m}(\text{e}[\text{new C}(\overline{v})/\text{this}][\overline{\text{u}}/\overline{\text{x}}])], \epsilon)}$$

$$return(\text{E}[\text{C.m}(v)]) = (\text{E}[v], \epsilon) \qquad \frac{\text{C.m} \in LibMeths \qquad \text{new C}(\overline{v}_1).\text{m}(\overline{\text{u}}_1) \xrightarrow{\epsilon}_{*} v}{call_{\text{C.m}}(\text{E}[\text{new C}(\overline{v}_1).\text{m}(\overline{\text{u}}_1)]) = (\text{E}[v], \epsilon)} \qquad \frac{\text{u.check}(v) \rightarrow^* v}{check(\text{E}[\text{u.check}(v)]) = (\text{E}[v], iev(v))}$$

$$\frac{\text{u.endorse}(v) \rightarrow^* v'}{endorse(\text{E}[\text{u.endorse}(v)]) = (\text{E}[v'], eev(v'))} \qquad if_{\mathbf{T}}(\text{E}[\text{if } v \text{ then } e_1 \text{ else } e_2]) = (\text{E}[e_1], \epsilon)$$

$$if_{\mathbf{F}}(\text{E}[\text{if } v \text{ then } e_1 \text{ else } e_2]) = (\text{E}[e_2], \epsilon) \qquad sequence(\text{E}[v; e]) = (\text{E}[e], \epsilon)$$

Fig. 17. Fundamental State Transformers Extracted from $\xrightarrow{\alpha}$.

$$com(\text{E}[\text{new C}(\overline{v}).\text{f}]) = select_{\text{f}} \qquad com(\text{E}[\text{new C}(\overline{v}).\text{m}(\overline{\text{u}})]) = call_{\text{C.m}} \quad \text{m} \notin \{\text{check}, \text{endorse}\}$$

$$com(\text{E}[\text{u.check}(v)]) = check \qquad com(\text{E}[\text{u.endorse}(v)]) = endorse \qquad com(\text{E}[\text{C.m}(v)]) = return$$

$$com(\text{E}[\text{if T then } e_1 \text{ else } e_2]) = if_{\mathbf{T}} \qquad com(\text{E}[\text{if F then } e_1 \text{ else } e_2]) = if_{\mathbf{F}}$$

Fig. 18. Definition of *com* for FJ$_{\text{taint}}$.

actual **T** or **F** value of the guard, and eliminating the effects of indirect flow from state transformation functions.

**Definition 5.6.** *The state transformers for FJ*$_{\text{taint}}$ *are composed of commands of the form select*$_{\text{f}}$ *for all fields f (selection), call*$_{\text{C.m}}$ *for all class, method pairs* C.m *(method dispatch), return (method return), endorse (endorsement), check (successful taint check within an sso), sequence (sequencing), and if*$_{\mathbf{T}}$ *and if*$_{\mathbf{F}}$ *(branch inlining). The behavior of these fundamental extracted state transformers are defined in Figure 17.*

Our treatment of library methods, check, and endorse bear discussion since they consider these in an atomic, "big step" manner. As noted in Section 3.2.2, when taint is propagated by library methods, for efficiency or implementation convenience it may be the case that taint propagation is not correctly applied until computed results are returned. This includes check and endorse, since technically these are library methods as per the definition in Section 3.1.6. Thus we specify that the extracted state transformer of any library method treat it atomically with respect to internal computations. In addition to *check* and *endorse*, for library methods where no security related events can occur we define a class of state transformers *call*$_{\text{C.m}}$ for C.m $\in$ *LibMeths*. This definition will also significantly simplify *our* proofs, and is irrelevant from a formal perspective since this definition yields the same observable events that a strict "small-step" definition of state transformers would for a given top-level program in the image of *Phos*.

Next, we define *com*, *state*, and $\langle \cdot, \cdot \rangle$ for FJ$_{\text{taint}}$. The command associated with a particular configuration e can be determined from its redex. We take the state of a configuration e to just be e itself, and combining a command and a state to obtain a configuration requires that the given command matches the form of the redex in the state– i.e. compatibility of the command and the state.

$\langle select_f, \mathrm{E}[\texttt{new}\ \mathrm{C}(\overline{\mathrm{v}}).\mathrm{f}]\rangle = \mathrm{E}[\texttt{new}\ \mathrm{C}(\overline{\mathrm{v}}).\mathrm{f}]$      $\langle call_{\mathrm{C.m}}, \mathrm{E}[\texttt{new}\ \mathrm{C}(\overline{\mathrm{v}}).\mathrm{m}(\overline{\mathrm{u}})]\rangle = \mathrm{E}[\texttt{new}\ \mathrm{C}(\overline{\mathrm{v}}).\mathrm{m}(\overline{\mathrm{u}})]$      $\langle return, \mathrm{E}[\mathrm{C.m}(\mathrm{v})]\rangle = \mathrm{E}[\mathrm{C.m}(\mathrm{v})]$

$\langle check, \mathrm{E}[\mathrm{u.check}(\mathrm{v})]\rangle = \mathrm{E}[\mathrm{u.check}(\mathrm{v})]$      $\langle endorse, \mathrm{E}[\mathrm{u.endorse}(\mathrm{v})]\rangle = \mathrm{E}[\mathrm{u.endorse}(\mathrm{v})]$

$\langle if_{\mathbf{T}}, \mathrm{E}[\texttt{if}\ \mathrm{v}\ \texttt{then}\ \mathrm{e}_1\ \texttt{else}\ \mathrm{e}_2]\rangle = \mathrm{E}[\texttt{if}\ \mathrm{v}\ \texttt{then}\ \mathrm{e}_1\ \texttt{else}\ \mathrm{e}_2]$      $\langle if_{\mathbf{F}}, \mathrm{E}[\texttt{if}\ \mathrm{v}\ \texttt{then}\ \mathrm{e}_1\ \texttt{else}\ \mathrm{e}_2]\rangle = \mathrm{E}[\texttt{if}\ \mathrm{v}\ \texttt{then}\ \mathrm{e}_1\ \texttt{else}\ \mathrm{e}_2]$

Fig. 19. Definition of $\langle \cdot, \cdot \rangle$ for FJ$_{\text{taint}}$

$$x =_{\circledast} x \qquad v =_{\circ} v \qquad v_1 =_{*} v_2 \qquad Op(\overline{e}) =_{\circledast} Op(\overline{e}) \qquad \frac{e_1 =_{\circledast} e_2}{e_1.f =_{\circledast} e_2.f} \qquad \frac{e_1, \overline{e}_1 =_{\circledast} e_2, \overline{e}_2}{e_1.m(\overline{e}_1) =_{\circledast} e_2.m(\overline{e}_2)}$$

$$\frac{e_1 =_{\circledast} e_2}{C.m(e_1) =_{\circledast} C.m(e_2)} \qquad \frac{\overline{e}_1 =_{\circledast} \overline{e}_2}{\texttt{new}\ C(t, \overline{e}_1) =_{\circledast} \texttt{new}\ C(t, \overline{e}_2)} \qquad \frac{C \in BaseTypes}{\texttt{new}\ C(\bullet, v_1) =_{\circ} \texttt{new}\ C(\bullet, v_2)} \qquad \frac{e_1 =_{\circledast} e_1' \qquad e_2 =_{\circledast} e_2'}{e_1; e_2 =_{\circledast} e_1'; e_2'}$$

$$\frac{e_1^1 =_{\circledast} e_1^2 \quad \cdots \quad e_n^1 =_{\circledast} e_n^2}{e_1^1 \cdots e_n^1 =_{\circledast} e_1^2 \cdots e_n^2} \qquad \frac{e_1 =_{\circledast} e_1' \qquad e_2 =_{\circledast} e_2' \qquad e_3 =_{\circledast} e_3'}{\texttt{if}\ e_1\ \texttt{then}\ e_2\ \texttt{else}\ e_3 =_{\circledast} \texttt{if}\ e_1'\ \texttt{then}\ e_2'\ \texttt{else}\ e_3'}$$

Fig. 20. Definition of Trust Equivalence and Shape Conformance Relations on Expressions

**Definition 5.7.** *We define* $state(\mathrm{e}) = \mathrm{e}$ *and define com as in Figure 18. We define* $\langle \cdot, \cdot \rangle$ *as in Figure 19.*

These definitions clearly satisfy the model requirements.

**Lemma 5.1.** *For any FJ*$_{\text{taint}}$ *configuration* $\kappa$, *we have* $\langle com(\kappa), state(\kappa) \rangle = \kappa$.

*Trust equivalence and state transformation.* Now we define trust equivalence $=_{\circ}$ on FJ$_{\text{taint}}$ expressions as required. This definition requires structural conformance of related states (expressions), and requires agreement of base values except in the case of tainted base objects. Aside from satisfying the model definition, the definition of trust equivalence will be crucial in our proof of explicit integrity modulo endorsement, as it defines the the necessary inductive invariant on extracted state transformations for this result.

Also, since endorsement may allow trust equivalent states to transform into non-structural equivalence, to satisfy the completeness requirement of Definition 5.3 we need to show that transformation preserves a weaker structural conformance relation $=_{*}$ on states (expressions). These relations are very similar with $=_{*}$ strictly weaker than $=_{\circ}$, and in proofs we will generally consider them together. Hence we define the metavariable $=_{\circledast}$ to range over $=_{*}$ and $=_{\circ}$.

**Definition 5.8.** *The* trust equivalence $=_{\circ}$ *and* shape conformance $=_{*}$ *relations on expressions are defined as the least relations inductively satisfying the rules in Figure 20, where* $=_{\circledast}$ *is a metavariable that ranges over* $=_{\circ}$ *and* $=_{*}$.

### 5.2.1. Sanity Conditions on Library Methods

We define two sanity conditions for library methods: *not undertainting* and *not overtainting*. The former condition is required in the implementation in order to meet explicit integrity modulo endorsement,

whereas the latter is a good practice in the implementation of taint analysis tools. Hereafter we will assume that library methods are not undertainting.

**Definition 5.9.** *We say* C.m $\in$ *LibMeths is not undertainting iff for all* $\overline{v}_1$, $\overline{u}_1$, $\overline{v}_2$, $\overline{u}_2$ *where:*

$$\overline{v}_1, \overline{u}_1 =_o \overline{v}_2, \overline{u}_2 \qquad call_{\text{C.m}}(\text{new } \text{C}(\overline{v}_1).\text{m}(\overline{u}_1)) = (v_1, \epsilon) \qquad call_{\text{C.m}}(\text{new } \text{C}(\overline{v}_2).\text{m}(\overline{u}_2)) = (v_2, \epsilon)$$

*we have* $v_1 =_o v_2$.

For example, String.concat is not undertainting if the taint propagation policy is defined as in Section 3.2.2 where the taint of a concatenated string is the meet of its operands' taints, but it would be e.g. if its results were always untainted.

Not overtainting refines the precision of taint tracking with respect to a given state. Intuitively, a library method that only *directly* depends on its high integrity inputs is not overtainting if its results are untainted.

**Definition 5.10.** *We say* C.m $\in$ *LibMeths is not overtainting with respect to input* $\overline{v}_1, \overline{u}_1$ *iff for all* $\overline{v}_2$, $\overline{u}_2$ *where:*

$$\overline{v}_1, \overline{u}_1 =_o \overline{v}_2, \overline{u}_2 \qquad call_{\text{C.m}}(\text{new } \text{C}(\overline{v}_1).\text{m}(\overline{u}_1)) = (v_1, \epsilon) \qquad call_{\text{C.m}}(\text{new } \text{C}(\overline{v}_2).\text{m}(\overline{u}_2)) = (v_2, \epsilon)$$

*if* $v_1 = v_2$ *then* $v_1 = \text{new } \text{D}(o, \overline{v})$ *for some* $\overline{v}$.

### 5.3. Extended Example

Assume given a class table *CT* containing sanitizer and sso methods Sec.sanitize and Sec.secureMeth which are identity function for the sake of brevity, i.e.:

$$mbody_{CT}(\text{Sec}, \text{sanitize}) = \text{x}, \text{x} \qquad mbody_{CT}(\text{Sec}, \text{secureMeth}) = \text{x}, \text{x}$$

and let $mbody_{CT}(\text{main}, \text{TopLevel}) = \text{attack}, \text{e}$ where e is:

```
if attack.eq(new String("foo")) then
    new Sec().secureMeth(attack)
else
    new Sec().secureMeth(new Sec().sanitize(attack))
```

Note that this is an example of a program that is unsafe by our definition, since a tainted value can flow directly into an sso, though it is noninterfering modulo endorsement since that value can only be new String("foo") (similar to the example at the beginning of this Section). However, *Phos* will emplace a check that will ensure blocking of unsafe executions. We note that the execution of *Phos*($\mathfrak{p}$(new String("foo"))) up to the point it gets stuck within Sec.check is associated with the following state transformer *f*:

$$f = sequence * return * call_{\text{Sec.log}} * call_{\text{Sec.secureMeth}} * if_{\textbf{T}} * call_{\text{String.eq}}$$

Observe that $\pi_2(f(Phos(\mathfrak{p}(\mathbf{a})))) = \epsilon$ for any $\mathbf{a}$, trivially satisfying the requirements of explicit integrity modulo endorsement. Crucially, note that $\mathbf{a}$ need not be the string $''\texttt{foo}''$ in order for $f(Phos(\mathfrak{p}(\mathbf{a})))$ to be defined– even though the program $Phos(\mathfrak{p}(\mathbf{a}))$ would not take the $\mathbf{T}$ branch through the conditional during actual execution, it is "forced" that way by $f$. This is central to the definition of explicit attacker power with respect to $Phos(\mathfrak{p}(\texttt{new String}(''\texttt{foo}'')))$ and $f$.

In contrast, the state transformer associated with the actual execution of $Phos(\mathfrak{p}(\texttt{new String}(s)))$ for $s \neq ''\texttt{foo}''$ up to the point it gets endorsed by $\texttt{String.endorse}$ within $\texttt{Sec.sanitize}$ is:

$$g = endorse * call_{\texttt{Sec.sanitize}} * if_{\mathbf{F}} * call_{\texttt{String.eq}}$$

We note that:

$$\pi_2(g(Phos(\mathfrak{p}(\texttt{new String}(s'))))) = eev(\texttt{new String}(s'))$$

for all $s'$. Furthermore, continued execution of $Phos(\mathfrak{p}(\texttt{new String}(s)))$ is associated with the following function $h$ which takes the program through the successful check of the sanitized object:

$$h = check * sequence * return * call_{\texttt{Sec.log}} * call_{\texttt{Sec.secureMeth}}$$

We note that:

$$\pi_2(h * g(Phos(\mathfrak{p}(\texttt{new String}(s'))))) = eev(\texttt{new String}(s')), iev(\texttt{new String}(s'))$$

for all $s'$. Finally, we observe that:

$$p_e(Phos(\mathfrak{p}(\mathbf{a})), g) = p_e(Phos(\mathfrak{p}(\mathbf{a})), h * g) = \{Phos(\mathfrak{p}(\mathbf{a}))\}$$

for all $\mathbf{a}$, satisfying the requirements of explicit integrity modulo endorsement.

Since $f$ and $h * g$ represent all possible control flow paths through the program that can generate events, it is evident that $Phos(\mathfrak{p}(\mathbf{a}))$ satisfies explicit integrity modulo endorsement for all $\mathbf{a}$.

### 5.4. Enforcement of Explicit Integrity Modulo Endorsement by Phos

Our general strategy is to show that non-endorsement events do not change attacker power as required by definition of explicit integrity modulo endorsement. We begin with useful substitution Lemmas for expressions and contexts showing that substitution in terms and in contexts preserves shape conformance and trust equivalence.

**Lemma 5.2.** *If* $v_1 =_\circledast v_2$ *then* $e[v_1/x] =_\circledast e[v_2/x]$[5].

**Lemma 5.3.** *If* $E_1[e_1] =_\circledast E_2[e_2]$ *and* $e_1$, $e_2$ *are both redexes then* $e_1 =_\circledast e_2$, *and if* $e_1' =_\circledast e_2'$ *then* $E_1[e_1'] =_\circledast E_2[e_2']$.

---

[5]Since $=_\circledast$ is a metavariable, this and subsequent results hold for consistent substitution of either $=_\circ$ or $=_*$ for $=_\circledast$.

Now, we show that both shape conformance and trust equivalence is preserved by all non-endorsement events. Note that since both $=_\circ$ and $=_*$ are symmetric it suffices to prove this and subsequent results as implications. Since we assume that library methods are defined on all well-formed inputs and that their results preserve conformance and trust by virtue of sanity conditions, it is adequate to not explicitly consider transformations *within* library methods, hence we define $\dot{f}$ iff for all $e$ there does not exist $E$, $e$, and $\alpha$ such that $f(e) = (C.m(e), \alpha)$.

**Lemma 5.4.** *If* $e_1 =_\circledast e_2$ *and* $\dot{f}(e_1) = (e_1', \alpha_1)$ *where* $\alpha_1 \notin \mathbf{E_{en}}^*$, *then* $\dot{f}(e_2) = (e_2', \alpha_2)$ *where:*

(1) $e_1' =_\circledast e_2'$
(2) *if* $e_1 =_\circ e_2$ *then* $\alpha_1 = \alpha_2$.

**Proof.** If suffices to consider just the cases where $\dot{f}$ is a fundamental state transformer as defined in Figure 17, since any state transformer will be composed of these.

*Case* $select_{f_i}$. In this case by definition of $\langle \cdot, \cdot \rangle$ and $=_\circledast$ we have that $e_1$ and $e_2$ are of the form $E_1[\texttt{new } C(t, \overline{v}).\texttt{f}]$ and $E_2[\texttt{new } C(t, \overline{u}).\texttt{f}]$ respectively, and since primitive fields can only be selected within library methods therefore $C$ is not a base type, thus $\overline{v} =_\circledast \overline{u}$ by Lemma 5.3 and definition of $=_\circledast$ (either as $=_\circ$ or $=_*$). But then $e_1' = E_1[v_i]$ for $v_i \in \overline{v}$ and $\alpha_1 = \epsilon$ by definition of $select_{f_i}$, and letting $e_2' = E_2[u_i]$ for $u_i \in \overline{u}$ we have $select_{f_i}(e_2) = (e_2', \epsilon)$ by definition. And $v_i =_\circledast u_i$ in this case by definition of $=_\circledast$, hence $E_1[v_i] =_\circledast E_2[u_i]$ by Lemma 5.3, so the result follows in this case.

*Case* $call_{C.m}$. In this case by definition of $\langle \cdot, \cdot \rangle$ and $=_\circledast$ we have that $e_1$ and $e_2$ are of the form $E_1[\texttt{new } C(\overline{v}_1).\texttt{m}(\overline{v}_2)]$ and $E_2[\texttt{new } C(\overline{u}_1).\texttt{m}(\overline{u}_2)]$ respectively, where $\overline{v}_1 =_\circledast \overline{u}_1$ and $\overline{v}_2 =_\circledast \overline{u}_2$ by Lemma 5.3. But $C.m \notin \textit{LibMeths}$ by assumption, so by definition of $call_{C.m}$ where $mbody_{CT}(C, m) = \overline{x}, e$ we have:

$$e_1' = E_1[e[\texttt{new } C(\overline{v}_1)/\texttt{this}][\overline{v}_2/\overline{x}]] \qquad \alpha_1 = \epsilon$$

and letting:

$$e_2' = E_2[e[\texttt{new } C(\overline{u}_1)/\texttt{this}][\overline{u}_2/\overline{x}]] \qquad \alpha_2 = \epsilon$$

we have $call_{C.m}(e_2) = (e_2', \alpha_2)$, therefore the result follows by Lemmas 5.2 and 5.3 since $\overline{v}_2 =_\circledast \overline{u}_2$ as reasoned above.

*Case check.* In this case by definition of $\langle \cdot, \cdot \rangle$ and $=_\circledast$ we have that $e_1$ and $e_2$ are of the form $E_1[u_1.\texttt{check}(v_1)]$ and $E_2[u_2.\texttt{check}(v_2)]$ where $v_1 =_\circledast v_2$ by definition of $=_\circledast$ and Lemma 5.3. Now, by definition of *check* we have that $e_1' = E_1[v_1]$, and $v_1$ must be an untainted or maybe tainted value otherwise the implicit call to check would block, and also $\alpha_1 = iev(v_1)$. But $v_2$ is at the same taint level as $v_1$ by definition of $=_\circledast$, so letting $e_2' = E_2[v_2]$ and $\alpha_2 = iev(v_2)$ we have $check(e_2) = (e_2', \alpha_2)$ by definition, and $e_1' =_\circledast e_2'$ by Lemmas 5.2 and 5.3. Finally, since both $v_1$ and $v_2$ are not tainted, if $=_\circledast$ is $=_\circ$ then they must be equivalent in which case $\alpha_1 = \alpha_2$.

*Case* $if_{\mathbf{T}}$. In this case by definition of $\langle \cdot, \cdot \rangle$ and $=_\circledast$ we have that $e_1$ and $e_2$ are of the form $E_1[\texttt{if } v \texttt{ then } e_1 \texttt{ else } ...]$ and $E_2[\texttt{if } u \texttt{ then } e_2 \texttt{ else } ...]$ where $e_1 =_\circledast e_2$ by definition of $=_\circledast$ and Lemma 5.3. By definition of $if_{\mathbf{T}}$ we have that $e_1' = E_1[e_1]$ and $\alpha_1 = \epsilon$, and letting $e_2' = E_2[e_2]$ and $\alpha_2 = \epsilon$ we have $if_{\mathbf{T}}(e) = (e_2', \alpha_2)$. Hence the result follows in this case by 5.3. $\square$

Next, to generalize to all transformers including endorsement, we show that shape conformance is preserved by all transformations, and also that if transformations of two trust equivalent states yields the same observable events, then the resulting states will also be trust equivalent.

**Lemma 5.5.** *If* $e_1 =_* e_2$ *and* $\dot{f}(e_1) = (e_1', \alpha_1)$ *then* $\dot{f}(e_2) = (e_2', \alpha_2)$ *where:*

(1)  $e_1' =_* e_2'$
(2)  *if* $e_1 =_\circ e_2$ *and* $\alpha_1 = \alpha_2$ *then* $e_1' =_\circ e_2'$.

**Proof.**  As for Lemma 5.4, it suffices to consider just the fundamental state transformers. Lemma 5.4 also covers all cases for this Lemma except *endorse*, which we demonstrate as follows.

*Case endorse.*   In this case by definition of $\langle \cdot, \cdot \rangle$ and $=_*$ we have that $e_1$ and $e_2$ are of the form $E_1[\texttt{new D}(t, \overline{v}_1).\texttt{endorse}()]$ and $E_2[\texttt{new D}(t, \overline{v}_2).\texttt{endorse}()]$ respectively, where at least $\overline{v}_1 =_* \overline{v}_2$ by definition. By definition of *endorse* we have:

$$e_1' = E_1[\texttt{new D}(t \vee \odot, \overline{v}_1)] \qquad\qquad \alpha_1 = eev(\texttt{new D}(t \vee \odot, \overline{v}_1))$$

and letting:

$$e_2' = E_2[\texttt{new D}(t \vee \odot, \overline{v}_2)] \qquad\qquad \alpha_2 = eev(\texttt{new D}(t \vee \odot, \overline{v}_2))$$

we have by definition that *endorse*$(e_2) = (e_2', \alpha_2)$ and also $e_1' =_* e_2'$ by Lemma 5.3. Furthermore, if we assume that $\alpha_1 = \alpha_2$, this implies that $\overline{v}_1 = \overline{v}_2$, hence if $e_1 =_\circ e_2$ then also $e_1 =_\circ e_2$ by Lemmas 5.2 and 5.3.                                                                     $\square$

The required completeness property of $FJ_{\text{taint}}$ state transformers now falls out as a corollary of Lemmas 5.4 and 5.5.

**Lemma 5.6.** *All FJ$_{\text{taint}}$ state transformers are complete in the sense of Definition 5.3.*

As an auxiliary Lemma to our main result, we show that attacker power is monotonic with respect to composition of state transformations.

**Lemma 5.7.** $p_e(e, f) \supseteq p_e(e, g * f)$ *for all f and g.*

As another auxiliary Lemma we show that when considering attacker power, it suffices to consider just transformations $\dot{f}$ that do not end in the middle of library methods as in the above Lemmas.

**Lemma 5.8.** *For all* $e$ *and* $f$, *there exists* $\dot{f}$ *such that* $p_e(e, f) = p_e(e, \dot{f})$.

Now we can leverage the above results to establish a simple proof by contradiction of our main Theorem for this Section.

**Theorem 5.1.** *If* $e$ *is in the image of Phos, then it enjoys explicit integrity modulo endorsement.*

**Proof.** Suppose on the contrary that $e$ did not enjoy explicit integrity modulo endorsement. Then by definition and Lemma 5.8 there exists $\dot{f}$ and $\dot{g}$ such that:

$$\dot{f}(e) = (e_1, \alpha_1) \qquad \dot{g}(e_1) = (e_2, \alpha_2) \qquad \alpha_2 \notin \mathbf{E_{en}}^* \qquad p_e(e, \dot{f}) \neq p_e(e, \dot{g} * \dot{f})$$

But the last inequality and Lemma 5.7 imply that $p_e(e, \dot{f}) \supset p_e(e, \dot{g} * \dot{f})$, hence there exists some $e'$ such that:

$$e =_\circ e' \qquad \dot{f}(e') = (e'_1, \alpha_1) \qquad \dot{g}(e'_1) = (e'_2, \alpha'_2) \qquad \alpha_2 \neq \alpha'_2$$

and by Lemma 5.5 it must be the case that $e_1 =_\circ e'_1$, so by Lemma 5.4 we have $\alpha_2 = \alpha'_2$, which is a contradiction. $\qquad\square$

## 6. An Implementation of *Phos* in OpenMRS

In Section 1.1 we discussed an XSS vulnerability in the OpenMRS system (corrected in the current version) that inspired our interest in an in-depth taint analysis to better track data flow into secure operations and to enforce some level of sanitization. To explore and evaluate our proposed methods in practice, we have developed an automated analysis for OpenMRS by direct modification of the Phosphor system [5]. Our modification supports dynamic integrity taint analysis both prospectively and retrospectively. Our implementation is based on the formal model developed in previous sections, which enjoys a correctness guarantee. In this Section we describe our implementation and our evaluation of it.

### 6.1. Modifications to Phosphor

Out of the box, Phosphor provides a binary taint labeling scheme, with no support for endorsement. Users specify their security policy by identifying high integrity sinks, which are then automatically instrumented at the bytecode level with checks for low integrity inputs, by a combination of program rewriting and runtime mechanisms. Thus, to implement our in-depth taint analysis specification we needed to generalize the taint labeling scheme, add an endorsement mechanism, and add support for audit logging to the existing Phosphor codebase. This yielded our *Phos* implementation, as distinct from Phosphor.

Phosphor distinguishes only between two types of data– tainted and untainted. To support a generalized labeling scheme, in *Phos* we added to the Phosphor `Taint` class definition a field containing a `TaintLevel` enumeration. This latter type is endowed with a partial ordering that is specified by the programmer via an underlying graph definition, and join and meet operations. In our implementation we support the taint label lattice defined in Section 3.2 but this could be easily changed to accommodate others. We also define an `endorse` operation that takes the join of the input taint label and `MaybeTainted` as in this paper. Since Phosphor itself adds a `Taint` object to all program objects, these modifications are propagated through the system by the existing codebase.

As for ordinary Phosphor, in *Phos* we allow specification of secure sinks, however the rewriting algorithm adds instrumentation for audit logging of values at or below a specified taint level that reach any sink (`MaybeTainted` in our case). The following information is logged in such a case: the function name of any sink that had a tainted variable pass through it, the taint level of any sunken tainted variable,
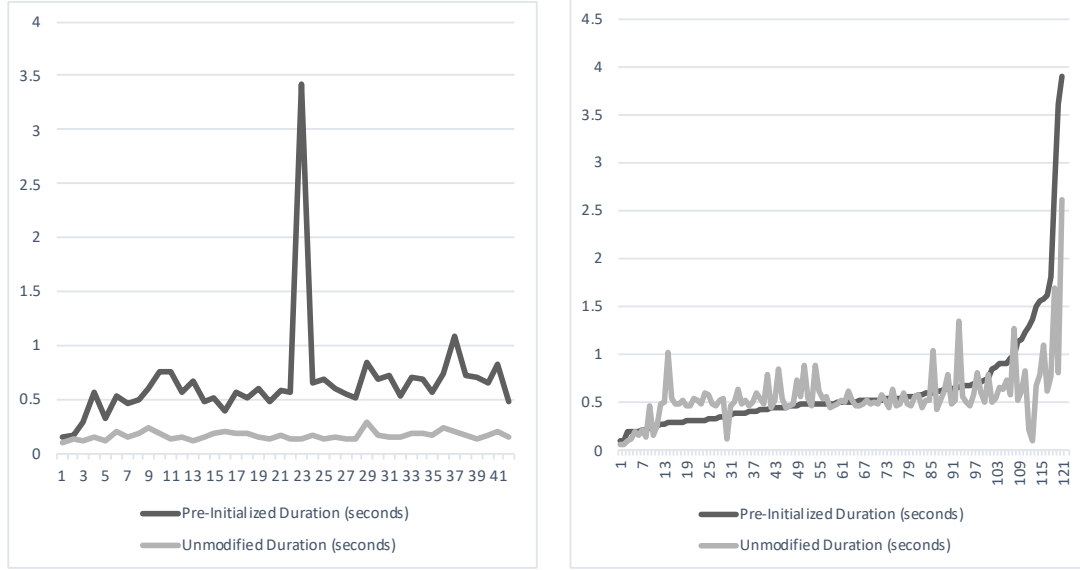
Fig. 21. *Phos* instrumentation timing overhead for OpenMRS for actions (left) and page loads (right). Numbers on the x-axis identify particular actions and page loads, and the y-axis is completion time in seconds.

the value of sunken variables, and A stack trace of the thread when a tainted variable was sunk. Much of this information was already being collected in the unmodified Phosphor.

We also allow specification of a set of sanitizers in the same manner as secure sink specification, i.e. specific methods are identified in an initial configuration file provided when rewriting a program. These methods will have return values `endorse`d via insertion of that method. Thus the end product functions the same as the system specified in Section 4– an input set of *SSOs* and *Sanitizers* are provided, along with a program for instrumentation, and the program is rewritten with instrumentation to support $SP_{taint}$ an $LS_{taint}$. Our *Phos* implementation also supports a specification of low integrity sources at arbitrary taint levels. The implementation is available on a public GitHub [50], as well as our *Phos*-instrumented version of OpenMRS.

### 6.2. OpenMRS Sources, Sinks, and Sanitizers

To apply *Phos* to OpenMRS, it is necessary to identify sources, sinks, and sanitizers in the system. Since our concern is mainly defense against injection and XSS attacks, we focused on database inter-actions. OpenMRS in its current form uses the popular Hibernate ORM framework as a database API, which supports two ways of interacting with databases– via persistent relationally mapped object saving/loading, and via queries. For our scope of work we focused on queries based on data in memory rather than persistent data, since the latter would require persistence of taint information and hence a far more complex implementation task.

The lists of sources, sinks, and sanitizers we identified in OpenMRS are provided in our implementation GitHub [50]. Our method for identifying sinks and sanitizers was to leverage our knowledge of the Hibernate API. Specifically, to identify sinks, we searched the OpenMRS codebase for methods that employ Hibernate database write functionality. To identify sanitizers, we searched the OpenMRS code-
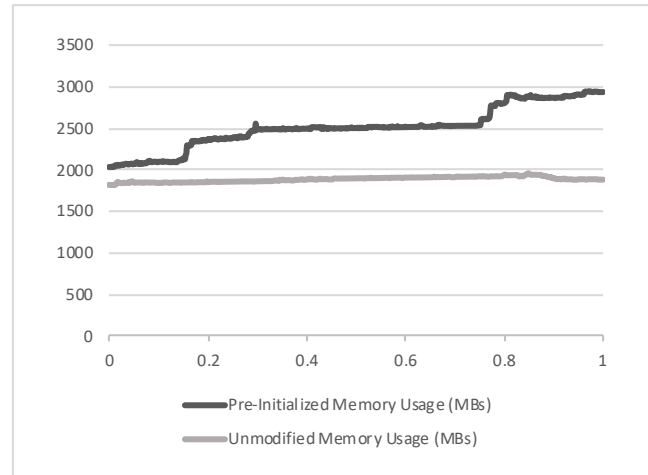
Fig. 22. Phosphor instrumentation memory overhead. The x-axis denotes the fraction completion of a test run, and the y-axis denotes memory used in MB.

base for methods that employ Hibernate sanitization functionality. The list of sources was determined by searching for methods that use `javax.servlet` functionality for recovering data from `POST` requests.

Another subtle but important detail of our integration of *Phos* with OpenMRS is that in OpenMRS, the arguments for the sinks are not necessarily tainted themselves, but rather are objects containing tainted member variables. Therefore, we also modified Phosphor to not only check sink arguments for a taint, but also argument member variables.

## 6.3. Implementation Evaluation

To evaluate our instrumentation of OpenMRS with *Phos*, we developed an automated testing method to evaluate correctness of the implementation, as well as timing and memory overhead. However, to understand our evaluation, certain details need to be explained.

*Phosphor Initialization Overhead.* When instrumenting Java classes with Phosphor, one can either run the software manually, specifying all of the source files to the program, or Phosphor can automatically detect and instrument uninstrumented code as the program runs. The latter option was chosen for this project due to the nature of OpenMRS and the onerous overhead of manual instrumentation.

As a consequence of instrumenting Java classes dynamically, an instrumentation overhead occurs as new uninstrumented source code is discovered. Thus, an initial run through a Phosphor-modified OpenMRS would be slower than consecutive runs, which was indeed observed by testing– the initial run of an particular method typically took about twice as long as subsequent runs.

Since our main concern in evaluation was to compare the overhead of instrumented vs. unmodified OpenMRS, and initialization overhead is arguably amortized to insignificance over long use sessions, the results we report here are only for pre-initialized testing runs. However we do note this additional overhead with the use of Phosphor's dynamic instrumentation feature.

*Actions and Page Loads.* The first category obtained from running the Fuzzer is the action duration data. This data comes from submitting a form containing arbitrary data and logging how long it takes for the subsequent page is loaded, indicating that the data was processed and potentially sent through a

sanitizer and/or a sink. After visiting every page the Fuzzer could find, around 50 forms were submitted and logged for each run.

While the focus of this analysis pertains to the results obtained from submitting data to OpenMRS, the overhead occurred by tracking tainted data can also be observed in the time it takes to traverse between different pages, not just when traversing from a submit button. As such, the Fuzzer was implemented to additionally record the time it takes to navigate between pages. These page loading results and the averages/ratios are thus presented for each of the three runs.

### 6.3.1. Experiments and Results

To evaluate *Phos*-instrumented OpenMRS, we developed a script that iterated over 42 actions, and over 121 page loads, recording timing and memory use, that we call a *test run*. We did a test run over unmodified OpenMRS to establish a baseline, and also did a test run over OpenMRS instrumented with our implementation of *Phos*. Finally, to evaluate how much our modifications impact Phosphor overhead, we did an actions-only test run over OpenMRS instrumented with pre-initialized unmodified Phosphor.

An initial concern of our evaluation was determining whether the system worked correctly, and whether data reaching sinks was maybe tainted, indicating sanitization, as well as being logged properly. We confirmed this, and did not discover instances of unsanitized data reaching sinks. Subsequently, we considered timing and memory consumption.

*Timing.* Our timing results are summarized in Table 1. Here we show the average time to complete each action and page load for the unmodified OpenMRS baseline, as well as OpenMRS instrumented with our implementation of *Phos*, and OpenMRS instrumented with pre-initialized unmodified Phosphor. These results demonstrate that instrumentation imposes a bit less than 3x overhead, while average times for completion are not onerous. Furthermore, our comparison of *Phos* and Phosphor shows that our modifications to did not add significant overhead to the taint analysis.

Table 1

Average timing and overhead for unmodified OpenMRS (baseline), versus instrumented with Phosphor and with *Phos*.

|  | **Actions** | | **Loads** | |
| --- | --- | --- | --- | --- |
|  | Avg (secs) | Overhead | Avg (secs) | Overhead |
| OpenMRS Baseline | .236 | - | .567 | - |
| OpenMRS + Phosphor | .614 | 261% | - | - |
| OpenMRS + *Phos* | .670 | 284% | .636 | 112% |

Figure 21 shows more detailed results, comparing times for the OpenMRS baseline and the *Phos*-instrumented version for each action (left graph) and page load (right graph). These results show that timing overhead is fairly consistent, albeit with some significant anamolies. In particular, overhead for the instrumented version spiked on the `dataExport.list` action, which is action number 23 in the graph. It is unclear what caused this anomaly, but appears to be an artifact of the Phosphor implementation (not our modifications).

*Memory.* Figure 22 shows baseline memory consumption during test runs of unmodified OpenMRS, versus OpenMRS instrumented with pre-initialized unmodified Phosphor. As these results demonstrate, while instrumentation does impose memory overhead, the impact on performance is not practically significant.

# 7. Conclusion

In this paper we considered integrity taint analysis in a pure object-oriented language model. Our security model accounts for sanitization methods that may be incomplete, a known problem in practice and one inspired by our study of the OpenMRS medical records software system. We proposed an in-depth security mechanism based on combining prospective measures (to support access control) and retrospective measures (to support auditing and accountability) that address incomplete sanitization. More precisely, we propose treating the results of sanitization as "partially" endorsed, or "maybe tainted", and allow maybe tainted values to be useds in security sensitive operations but record such events in the audit log.

We developed a uniform security policy of dynamic integrity taint analysis that specifies both prospective and retrospective measures, separate from code. The specification is defined in terms of a logical interpretation of program traces and leverages techniques from information algebra, allowing prospective and retrospective measures to be characterized in a uniform and integrated manner. Since the specification is defined separate from code, we use it to establish provable correctness conditions for a rewriting algorithm that instruments in-depth integrity taint analysis. A rewriting approach supports development of tools that can be applied to legacy code without modifying language implementations.

Although our specification of dynamic integrity taint analysis with endorsement establishes correctness conditions for implementations, it is still operational in nature. We therefore developed the hyperproperty of explicit integrity modulo endorsement to characterize the security property of integrity taint analysis in a non-operational manner. It is important to note that this formulation was not simply the dualization of previous formulations of explicit secrecy [10], since these formulations address only low-level code with unstructured data. We subsequently demonstrated that the image of our rewriting algorithm enjoys this security property.

Since our broader goal is to support well-founded practical tools for hardening software, we developed an instrumented version of OpenMRS that integrates our in-depth taint analysis formally specified in our model. Results from our evaluation of this implementation suggest that it is correct and practically feasible. We have made the implementation available on a public GitHub [50].

# References

[1] Edward J. Schwartz, Thanassis Avgerinos, and David Brumley. All you ever wanted to know about dynamic taint analysis and forward symbolic execution (but might have been afraid to ask). In *IEEE S&P*, pages 317–331, 2010.

[2] Benjamin Livshits, Michael Martin, and Monica S Lam. Securifly: Runtime protection and recovery from web application vulnerabilities. Technical report, Technical report, Stanford University, 2006.

[3] Gary Wassermann and Zhendong Su. Sound and precise analysis of web applications for injection vulnerabilities. In *PLDI*, pages 32–41, 2007.

[4] Vinod Ganapathy, Trent Jaeger, Christian Skalka, and Gang Tan. Assurance for defense in depth via retrofitting. In *LAW*, 2014.

[5] Jonathan Bell and Gail E. Kaiser. Phosphor: illuminating dynamic data flow in commodity jvms. In *OOPSLA*, pages 83–101, 2014.

[6] Jonathan Bell and Gail E. Kaiser. Dynamic taint tracking for java with phosphor (demo). In *ISSTA*, pages 409–413, 2015.

[7] Atsushi Igarashi, Benjamin C. Pierce, and Philip Wadler. Featherweight java: a minimal core calculus for java and GJ. *ACM Trans. Program. Lang. Syst.*, 23(3):396–450, 2001.

[8] Sepehr Amir-Mohammadian, Stephen Chong, and Christian Skalka. Correct audit logging: Theory and practice. In *POST*, pages 139–162, 2016.

[9] Benjamin Livshits. Dynamic taint tracking in managed runtimes. Technical report, Technical Report MSR-TR-2012-114, Microsoft Research, 2012.

[10] Daniel Schoepe, Musard Balliu, Benjamin C. Pierce, and Andrei Sabelfeld. Explicit secrecy: A policy for taint tracking. In *IEEE EuroS&P*, pages 15–30, 2016.

[11] Michael R. Clarkson and Fred B. Schneider. Hyperproperties. *Journal of Computer Security*, 18(6):1157–1210, 2010.

[12] J. Kohlas. *Information Algebras: Generic Structures For Inference*. Discrete mathematics and theoretical computer science. Springer, 2003.

[13] Juerg Kohlas and Juerg Schmid. An algebraic theory of information: An introduction and survey. *Information*, 5(2):219–254, 2014.

[14] Sepehr Amir-Mohammadian and Christian Skalka. In-depth enforcement of dynamic integrity taint analysis. In *PLAS*, 2016.

[15] Michael Martin, Benjamin Livshits, and Monica S. Lam. Finding application errors using PQL: A program query language. In *OOPSLA*, 2005.

[16] R. Sekar. An efficient black-box technique for defeating web application attacks. In *NDSS*, 2009.

[17] Zheng Wei and David Lie. Lazytainter: Memory-efficient taint tracking in managed runtimes. In *SPSM Workshop at CCS*, pages 27–38, 2014.

[18] William Enck, Peter Gilbert, Byung-Gon Chun, Landon P. Cox, Jaeyeon Jung, Patrick McDaniel, and Anmol Sheth. Taintdroid: an information flow tracking system for real-time privacy monitoring on smartphones. *Commun. ACM*, 57(3):99–106, 2014.

[19] Prateek Saxena, R. Sekar, and Varun Puranik. Efficient fine-grained binary instrumentationwith applications to taint-tracking. In *CGO*, pages 74–83, 2008.

[20] Winnie Cheng, Qin Zhao, Bei Yu, and Scott Hiroshige. Tainttrace: Efficient flow tracing with dynamic binary rewriting. In *IEEE ISCC*, pages 749–754, 2006.

[21] Erik Bosman, Asia Slowinska, and Herbert Bos. Minemu: The world's fastest taint tracker. In *RAID*, pages 1–20, 2011.

[22] David (Yu) Zhu, Jaeyeon Jung, Dawn Song, Tadayoshi Kohno, and David Wetherall. Tainteraser: protecting sensitive data leaks using application-level taint tracking. *Operating Systems Review*, 45(1):142–154, 2011.

[23] Erika Chin and David Wagner. Efficient character-level taint tracking for java. In *ACM SWS*, pages 3–12, 2009.

[24] Vivek Haldar, Deepak Chandra, and Michael Franz. Dynamic taint propagation for java. In *ACSAC*, pages 303–311, 2005.

[25] Chiara Bodei and Letterio Galletta. Tracking sensitive and untrustworthy data in IoT. In *ITASEC*, pages 38–52, 2017.

[26] Dorothy E Denning and Peter J Denning. Certification of programs for secure information flow. *Communications of the ACM*, 20(7):504–513, 1977.

[27] Joseph A. Goguen and José Meseguer. Security policies and security models. In *IEEE S&P*, pages 11–20, 1982.

[28] Andrei Sabelfeld and Andrew C Myers. Language-based information-flow security. *IEEE Journal on selected areas in communications*, 21(1):5–19, 2003.

[29] Benjamin Livshits and Stephen Chong. Towards fully automatic placement of security sanitizers and declassifiers. In *POPL*, pages 385–398, 2013.

[30] Aslan Askarov and Andrei Sabelfeld. Gradual release: Unifying declassification, encryption and key release policies. In *IEEE S&P*, pages 207–221, 2007.

[31] Andrei Sabelfeld and David Sands. Declassification: Dimensions and principles. *Journal of Computer Security*, 17(5):517–548, 2009.

[32] Aslan Askarov, Sebastian Hunt, Andrei Sabelfeld, and David Sands. Termination-insensitive noninterference leaks more than just a bit. In *ESORICS*, pages 333–348, 2008.

[33] Aslan Askarov and Andrei Sabelfeld. Tight enforcement of information-release policies for dynamic languages. In *CSF*, pages 43–59, 2009.

[34] Aslan Askarov and Andrew Myers. A semantic framework for declassification and endorsement. In *ESOP*, pages 64–84, 2010.

[35] Dennis M. Volpano. Safety versus secrecy. In *SAS*, pages 303–311, 1999.

[36] Daniel Schoepe, Musard Balliu, Frank Piessens, and Andrei Sabelfeld. Let's face it: Faceted values for taint tracking. In *European Symposium on Research in Computer Security*, pages 561–580, 2016.

[37] Musard Balliu, Daniel Schoepe, and Andrei Sabelfeld. We are family: Relating information-flow trackers. In *European Symposium on Research in Computer Security*, pages 124–145, 2017.

[38] Andrew C. Myers, Andrei Sabelfeld, and Steve Zdancewic. Enforcing robust declassification and qualified robustness. *Journal of Computer Security*, 14(2):157–196, 2006.

[39] Arnar Birgisson, Alejandro Russo, and Andrei Sabelfeld. Unifying facets of information integrity. In *ICISS*, pages 48–65, 2010.

[40] Wilmer Ricciotti and James Cheney. Strongly normalizing audited computation. *CoRR*, abs/1706.03711, 2017.

[41] Francisco Bavera and Eduardo Bonelli. Justification logic and audited computation. *Journal of Logic and Computation*, page exv037, 2015.

[42] Fred B. Schneider. Enforceable security policies. *ACM Transactions on Information and System Security*, 3(1):30–50, 2000.

[43] Juerg Kohlas and Juerg Schmid. An algebraic theory of information: An introduction and survey. *Information*, 5(2):219–254, 2014.

[44] Usage statistics module. https://wiki.openmrs.org/display/docs/Usage+Statistics+Module, 2010. Accessed: 2015-09-27.

[45] Sepehr Amir-Mohammadian. *A Formal Approach to Combining Prospective and Retrospective Security*. PhD thesis, The University of Vermont, 2017.

[46] Lujo Bauer, Jarred Ligatti, and David Walker. More enforceable security policies. Technical Report TR-649-02, Princeton University, June 2002.

[47] Deepak Garg, Limin Jia, and Anupam Datta. Policy auditing over incomplete logs: Theory, implementation and applications. In *CCS 2011*, pages 151–162, 2011.

[48] J. G. Cederquist, Ricardo Corin, M. A. C. Dekker, Sandro Etalle, J. I. den Hartog, and Gabriele Lenzini. Audit-based compliance control. *International Journal of Information Security*, 6(2-3):133–151, 2007.

[49] Ulf Nilsson and Jan Maluszyynski. Definite logic programs. In *Logic, Programming and Prolog*, chapter 2. 2000.

[50] Christian Skalka, Sepehr Amir-Mohammadian, and Samuel Clark. Retrospective Taint Analysis for OpenMRS. https://github.com/uvm-plaid/phosphor-mod, 2019.