

# Combining Symbolic and Numerical Domains for Information Leakage Analysis

Agostino Cortesi<sup>1</sup>, Pietro Ferrara<sup>2</sup>, and Raju Halder<sup>3</sup> Matteo Zanioli<sup>4</sup>

<sup>1</sup> Ca' Foscari University, Venice, Italy

<sup>2</sup> Julia srl, Verona, Italy

<sup>3</sup> Indian Institute of Technology Patna, India

<sup>4</sup> Alpenite srl, Venice, Italy

**Abstract.** We introduce an abstract domain for information-flow analysis of software. The proposal combines variable dependency analysis with numerical abstractions, yielding to accuracy and efficiency improvements. We apply the full power of the proposal to the case of database query languages as well. Finally, we present an implementation of the analysis, called *Sails*, as an instance of a generic static analyzer. Keeping the modular construction of the analysis, the tool allows one to tune the granularity of heap analysis and to choose the numerical domain involved in the reduced product. This way the user can tune the information leakage analysis at different levels of precision and efficiency.

## 1 Introduction

Protecting the confidentiality is a relevant problem when sensitive information flows through computing systems or transmits over public networks. Standard protection mechanisms, such as encryption, access control, etc. can suitably be applied at source level, but they are unable to protect the confidentiality once the information is released from the source and is allowed to flow through the computing systems.

The starting point of secure information flow analysis in software applications is the classification of program variables into different security levels. In the simplest case, two levels are commonly used: public (or low,  $L$ ) and secret (or high,  $H$ ). The main purpose is to prevent the leakage of sensitive information when flowing (implicitly or explicitly) from a high variable  $h$  to a lower one  $l$ . An explicit flow from  $h$  to  $l$  occurs when the content of  $h$  directly affects (e.g., through an assignment operator)  $l$ . On the other hand, an implicit flow from  $h$  to  $l$  occurs when the content of  $l$  gets affected indirectly (e.g., through a boolean condition in an `if` statement) by  $h$ , as stated in [17].

There is a widespread literature on methods and techniques for checking secure information flows in software. Generally, works on information flow fall into two categories: (i) dynamic, instrumentation based approaches (e.g., tainting), and (ii) static, language-based approaches (e.g., type systems). The dynamic approaches introduce significant run-time overhead [10, 33]. The static

approaches typically require some changes to the language and the run-time environment as well as non-trivial type annotations [38], making their adoption too expensive in practice.

Nevertheless, despite of these deep and extensive works, their practical applications have been relatively poor. Usually these approaches work on an ad-hoc programming language [4], and they do not support mainstream languages. This means that one should completely rewrite a program in order to apply them to some existing code.

Recently a new generic static analyzer (Sample<sup>5</sup>) based on the Abstract Interpretation theory has been developed and applied to many different contexts and analysis. Roughly, this analyzer splits and combines the abstraction of the heap and the approximation of other semantic information, *e.g.* string [12], type [19] abstractions.

In this paper<sup>6</sup>, we introduce a language-based information-flow analysis of imperative and database query languages based on the Abstract Interpretation framework, by combining symbolic and numerical domains; we present the tool Sails (Static Analysis of Information Leakage with Sample); finally, we show experimental results applying Sails on security benchmark programs.

In particular,

1. we represent variables' dependences in the form of propositional formula  $\psi = x \rightarrow y$ , where  $x, y$  are variables and value of  $y$  possibly depend on the value of  $x$ ; in order to detect possible information leakage, we check the satisfiability of  $\psi$  when assigning each variable the truth value corresponding to its sensitivity level;
2. we define abstract semantics of (i) imperative and (ii) database query languages in the domain of propositional formulae, by considering an over-approximation of variables' dependences at each program point;
3. we enhance the accuracy of the technique by analysing programs over numerical abstract domains, using reduced product of the symbolic propositional formulae domain and numerical abstract domains;
4. finally, we show encouraging experimental results on a set of security benchmarks using the tool Sails which is implemented based on our proposal.

The overall analysis combines a symbolic variable dependency analysis, based on the propositional formulae domain [11], and a variable value dependency analysis using numerical abstractions (*e.g.*, intervals or polyhedra). Unlike other works, our proposal provides an information flow analysis without any major constraint on the target language, since it tracks information flows between variables and heap locations over programs written in mainstream object-oriented languages like Java and Scala.

The rest of the paper is organized as follows. Section 2 introduces the dependency analysis through the propositional formulae domain. Section 3 combines the dependency analysis with numerical domains through a reduced

---

<sup>5</sup> <http://www.pm.inf.ethz.ch/research/semper/Sample>

<sup>6</sup> The paper is a revised and extended version of [25, 47, 48]

product. An extension to the case of database query languages is discussed in Section 4. Section 5 presents the main issues we solved in order to plug this information flow analysis into `Sample` while developing `Sails`. Section 6 presents the experimental results when applying `Sails` to a complex case study and to the `SecuriBench-micro` suite. Finally, Section 7 presents the related work and Section 8 concludes.

## 2 Dependency Analysis

This section formalizes the dependency analysis and proves its soundness following the abstract interpretation framework.

### 2.1 The Language

For the sake of simplicity, we consider a simple imperative language where programs consist of labeled commands (similar to [26]). The syntax is defined in Table 1<sup>7</sup>.

Table 1: Syntax of the language

Expressions			
$\text{exp}$	$\in$	$E$	
$\text{exp}$	$::=$	$n$	where $n \in \mathbb{N}$
		$v$	
		$\text{exp}_1 \oplus \text{exp}_2$	where $\oplus = \{+, -, *, /\}$
Conditions			
$b$	$\in$	$B$	
$b$	$::=$	<code>true</code>	
		<code>false</code>	
		$b_1 \otimes b_2$	where $\otimes = \{\vee, \wedge\}$
		$\neg b$	
		$\text{exp}_1 \oslash \text{exp}_2$	where $\oslash = \{\leq, >, =\}$
Labeled commands			
$\ell$	$\in$	$L$	set of labels
$c$	$\in$	$C$	
$c$	$::=$	${}^\ell \text{skip}$	
		${}^\ell v := \text{exp}$	
		<code>if <math>{}^\ell b</math> then <math>c_1</math> else <math>c_2</math> <math>{}^{\ell'}</math> endif</code>	
		$c_1; c_2$	
		<code>while <math>{}^\ell b</math> do <math>c</math> <math>{}^{\ell'}</math> done</code>	
$\mathcal{P}$	$::=$	$c^\ell$	program that ends with label $\ell$

<sup>7</sup> In the rest of the paper, we will omit the initial and final labels of statements when not required.

Table 2: Initial label function

$in[\ell \text{ skip}] \stackrel{\text{def}}{=} \ell$
$in[\ell v := \text{exp}] \stackrel{\text{def}}{=} \ell$
$in[\ell \text{b then } c_1 \text{ else } c_2 \ell' \text{ endif}] \stackrel{\text{def}}{=} \ell$
$in[c_1; c_2] \stackrel{\text{def}}{=} in[c_1]$
$in[\text{while } \ell \text{b do } c \ell' \text{ done}] \stackrel{\text{def}}{=} \ell$

Table 3: Final label function

$\mathcal{P} ::= c^\ell$	$fin[\mathcal{P}] \equiv \ell$
$c ::= \ell \text{ skip}$	$fin[c] \equiv fin[\mathcal{P}]$
$\ell v := \text{exp}$	$fin[\ell v := \text{exp}] \equiv fin[c]$
$\text{if } \ell \text{b then } c_1 \text{ else } c_2 \ell' \text{ endif}$	$fin[\text{if } \ell \text{b then } c_1 \text{ else } c_2 \ell' \text{ endif}] \equiv fin[c]$
	$fin[c_1] \equiv \ell'$
	$fin[c_2] \equiv \ell'$
$c_1; c_2$	$fin[c_1; c_2] \equiv fin[c]$
	$fin[c_1] \equiv in[c_2]$
	$fin[c_2] \equiv fin[c]$
$\text{while } \ell \text{b do } c \ell' \text{ done}$	$fin[\text{while } \ell \text{b do } c \ell' \text{ done}] \equiv fin[c]$
	$fin[c] \equiv \ell$

Let  $in : \mathbf{C} \rightarrow \mathbf{L}$  and  $fin : \mathbf{C} \rightarrow \mathbf{L}$  be two functions. By  $in[c]$  and  $fin[c]$  we denote the *initial* and *final label* of command  $c \in \mathbf{C}$  respectively. These two functions are formally defined in Table 2 and Table 3.

Each command corresponds to one or more *actions*. The set of actions, denoted by  $\mathbf{A}$ , consists of  $\{\ell \text{ skip}, \ell v := \text{exp}, \ell \text{b}, \ell \neg \text{b}, \ell \text{endif}, \ell \text{done}\}$ . Let  $a : \mathbf{C} \rightarrow \wp(\mathbf{A})$  be a function that, given a command, returns the set of actions involved in it. The function  $a$  for various commands is defined in Table 4.

Without loss of generality, we assume that the variables appearing in a program are implicitly declared. We denote by  $V(\mathcal{P})$  the set of variables in program  $\mathcal{P}$  and, similarly, by  $V(\text{exp})$  and  $V(\text{b})$  the variables contained in expression  $\text{exp}$  and condition  $\text{b}$  respectively. The definition of  $V$  is reported in Table 5.

## 2.2 The Concrete Domain

An *environment*  $\rho \in \mathcal{E}$  is a function  $\rho : \mathbf{V} \rightarrow \mathbb{N}$  which assigns a value to each variable. A *state*  $\sigma \in \Sigma = (\mathbf{L} \times \mathcal{E})$  is a pair  $\langle \ell, \rho \rangle$  where the program label  $\ell$  is the label of the action to be executed and the environment  $\rho$  defines the values of program variables at  $\ell$ .

We denote by  $E[\text{exp}]\sigma$  and  $B[\text{b}]\sigma$  the evaluation of expression  $\text{exp} \in \mathbf{E}$  and condition  $\text{b} \in \mathbf{B}$  respectively on the state  $\sigma$ . The details can be found in Tables 6 and 7 respectively.

Table 4: Action function

$a[\text{skip}] \stackrel{\text{def}}{=} \{\text{skip}\}$
$a[v := \text{exp}] \stackrel{\text{def}}{=} \{v := \text{exp}\}$
$a[\text{if } b \text{ then } c_1 \text{ else } c_2 \text{ endif}] \stackrel{\text{def}}{=} \{b, \text{not } b, \text{endif}\} \cup a[c_1] \cup a[c_2]$
$a[c_1; c_2] \stackrel{\text{def}}{=} a[c_1] \cup a[c_2]$
$a[\text{while } b \text{ do } c \text{ done}] \stackrel{\text{def}}{=} \{b, \text{not } b, \text{done}\} \cup a[c]$

Table 5: Variables functions

$V(n) \stackrel{\text{def}}{=} \emptyset$
$V(v) \stackrel{\text{def}}{=} \{v\}$
$V(\text{exp}_1 \oplus \text{exp}_2) \stackrel{\text{def}}{=} V(\text{exp}_1) \cup V(\text{exp}_2)$
$V(\text{true}) \stackrel{\text{def}}{=} \emptyset$
$V(\text{false}) \stackrel{\text{def}}{=} \emptyset$
$V(b_1 \otimes b_2) \stackrel{\text{def}}{=} V(b_1) \cup V(b_2)$
$V(\text{exp}_1 \circ \text{exp}_2) \stackrel{\text{def}}{=} V(\text{exp}_1) \cup V(\text{exp}_2)$
$V(\text{skip}) \stackrel{\text{def}}{=} \emptyset$
$V(v := \text{exp}) \stackrel{\text{def}}{=} \{v\} \cup V(\text{exp})$
$V(\text{if } b \text{ then } c_1 \text{ else } c_2 \text{ endif}) \stackrel{\text{def}}{=} V(b) \cup V(c_1) \cup V(c_2)$
$V(c_1; c_2) \stackrel{\text{def}}{=} V(c_1) \cup V(c_2)$
$V(\text{while } b \text{ do } c \text{ done}) \stackrel{\text{def}}{=} V(b) \cup V(c)$

Given a program  $\mathcal{P}$ , the set of possible initial and final states are defined as  $\text{In}[\mathcal{P}] \equiv \{\langle \text{in}[\mathcal{P}], \rho \rangle \mid \rho \in \mathcal{E}\}$  and  $\text{F}[\mathcal{P}] \equiv \{\langle \text{fin}[\mathcal{P}], \rho \rangle \mid \rho \in \mathcal{E}\}$ .

The *labeled transition semantics*  $T[\mathbf{c}]$  of a command  $\mathbf{c} \in \mathcal{P}$  is a set of transitions  $\langle \sigma_1, \mathbf{a}, \sigma_2 \rangle$  between a state  $\sigma_1$  and its next states  $\sigma_2$  by an action  $\mathbf{a} \in a(\mathbf{c})$ . The triple  $\langle \sigma_1, \mathbf{a}, \sigma_2 \rangle$  is also denoted by  $\sigma_1 \xrightarrow{\mathbf{a}} \sigma_2$ . The transition function  $T : \mathbf{C} \rightarrow \wp(\Sigma \times \mathbf{A} \times \Sigma)$  in Table 8 tracks all reachable states.

A *labeled transition system* is a tuple  $\langle \Sigma, I, F, \mathbf{A}, T \rangle$ , where  $\Sigma$  is the set of states,  $I \subseteq \Sigma$  is a nonempty set of initial states,  $F \subseteq \Sigma$  is a set of final states,  $\mathbf{A}$  is a nonempty set of actions, and  $T \in \wp(\Sigma \times \mathbf{A} \times \Sigma)$  is the labeled transition relation.

We define the *partial trace semantics* of a transition system, similarly to [26], as the set of all possible traces of elements in  $\Sigma$  (denoted by  $\Sigma^*$ ), recording the observation of executions starting from initial states and possibly reaching final states in finite time.

$$\Sigma^* \in \wp(\Sigma \times \mathbf{A} \times \Sigma)$$

$$\Sigma^* = \{\sigma_0 \xrightarrow{\mathbf{a}_0} \dots \xrightarrow{\mathbf{a}_{n-1}} \sigma_n \mid n \geq 1 \wedge \sigma_0 \in I \wedge \forall i \in [0, n-1] : \sigma_i \xrightarrow{\mathbf{a}_i} \sigma_{i+1} \in T^\ell\}$$

Let  $\pi_0, \pi_1 \in \Sigma^*$  be two partial traces. We define the following lattice operators:

Table 6: Evaluation of expressions

$E \in \mathbb{E} \rightarrow (\mathcal{E} \rightarrow \mathbb{N})$ $E[[n]]\rho \stackrel{\text{def}}{=} n$ $E[[v]]\rho \stackrel{\text{def}}{=} \rho(v)$ $E[[\text{exp}_1 \oplus \text{exp}_2]]\rho \equiv v_1 \oplus v_2 \text{ (such that } v_1 = E[[\text{exp}_1]]\rho \wedge v_2 = E[[\text{exp}_2]]\rho)$
---

Table 7: Evaluation of boolean conditions

$B \in \mathbb{B} \rightarrow (\mathcal{E} \rightarrow \{\text{true}, \text{false}\})$ $B[[\text{true}]]\rho \stackrel{\text{def}}{=} \text{true}$ $B[[\text{false}]]\rho \stackrel{\text{def}}{=} \text{false}$ $B[[b_1 \otimes b_2]]\rho \stackrel{\text{def}}{=} b_1 \otimes b_2 \text{ (such that } b_1 = B[[b_1]]\rho \wedge b_2 = B[[b_2]]\rho)$ $B[[\text{exp}_1 \odot \text{exp}_2]]\rho \stackrel{\text{def}}{=} \begin{array}{l} \text{true if } \exists v_1 = E[[\text{exp}_1]]\rho : v_2 = E[[\text{exp}_2]]\rho : v_1 \odot v_2 \\ \text{false if } \exists v_1 = E[[\text{exp}_1]]\rho : v_2 = E[[\text{exp}_2]]\rho : \text{not}(v_1 \odot v_2) \end{array}$
---

- $\pi_0 \leq \pi_1$  if and only if  $\pi_0$  is a subtrace of  $\pi_1$ ,
- $\pi_0 \wedge \pi_1 = \pi$  such that  $(\pi \leq \pi_1) \wedge (\pi \leq \pi_2)$  and  $(\forall \pi' : (\pi' \leq \pi_1) \wedge (\pi' \leq \pi_2)) . \pi' \leq \pi$ .

$\Sigma^*$  equipped with the order relation “ $\leq$ ” and the meet operator “ $\wedge$ ”, forms the meet semi lattice  $\langle \Sigma^*, \leq, \wedge \rangle$ .

This partial trace semantics can be expressed in a fixpoint form as well.

$$\begin{aligned} \Sigma^* &= \text{lfp}^{\subseteq} F : \\ F \in \Sigma^* &\rightarrow \Sigma^* \\ F(X) &\stackrel{\text{def}}{=} \{ \sigma \xrightarrow{a'} \sigma' \in \mathbb{T} \mid \sigma \in \mathbb{U} \} \cup \\ &\{ \sigma_0 \xrightarrow{a_0} \dots \xrightarrow{a_{n-2}} \sigma_{n-1} \xrightarrow{a_{n-1}} \sigma_n \mid \sigma_0 \xrightarrow{a_0} \dots \xrightarrow{a_{n-2}} \sigma_{n-1} \in X \wedge \sigma_{n-1} \xrightarrow{a_{n-1}} \sigma_n \in \mathbb{T} \} \end{aligned}$$

Let  $\langle \wp(\Sigma^*), \subseteq, \emptyset, \Sigma^*, \cap, \cup \rangle$  be a complete lattice of partial execution traces, where “ $\subseteq$ ” is the classical subset relation, “ $\cup$ ” is the set union and “ $\cap$ ” the set intersection.

### 2.3 Abstract Domain: Pos

Among all the abstract domains which are used in abstract interpretation of logic programs, **Pos** has received considerable attention [2, 11]. This domain is most commonly applied to the analysis of groundness dependencies for logic programs.

Let  $\bar{V} = \{\bar{x}, \bar{y}, \bar{z}, \dots\}$  be a countably infinite set of propositional variables and let  $FP(\bar{V})$  be the set of all finite subsets of variables of  $\bar{V}$ . The set of propositional formulae containing variables in  $\bar{V}$  and logical connectives in  $\Gamma \subseteq \{\wedge, \vee, \rightarrow, \neg\}$  is denoted by  $\Omega(\Gamma)$ . Similarly, given  $U \in FP(\bar{V})$ , the set of propositional formulae containing variables in  $U$  and connectives in  $\Gamma$  is denoted by  $\Omega_U(\Gamma)$ .

Table 8: The transition function

$T[\text{skip}^{\ell t}] = \{\langle \ell, \rho \rangle \xrightarrow{\ell \text{skip}^{\ell t}} \langle \text{fin}[\text{skip}^{\ell t}], \rho \rangle \mid \rho \in \mathcal{E}\}$
$T[\text{v} := \text{exp}^{\ell t}] = \{\langle \ell, \rho \rangle \xrightarrow{\ell v := \text{exp}^{\ell t}} \langle \text{fin}[\text{v} := \text{exp}^{\ell t}], \rho[v \leftarrow v] \rangle \mid \rho \in \mathcal{E} \wedge v = E[\text{exp}]\rho\}$
$T[\text{if } \ell \text{b then } c_1 \text{ else } c_2 \text{ endif}^{\ell t}] = T[c_1] \cup T[c_2] \cup$ $\{\langle \ell, \rho \rangle \xrightarrow{\ell \text{b}} \langle \text{in}[c_1], \rho \rangle \mid \rho \in \mathcal{E} \wedge \text{true} = B[\text{b}]\rho\} \cup$ $\{\langle \ell, \rho \rangle \xrightarrow{\ell \text{not b}} \langle \text{in}[c_2], \rho \rangle \mid \rho \in \mathcal{E} \wedge \text{false} = B[\text{b}]\rho\} \cup$ $\{\langle \ell', \rho \rangle \xrightarrow{\ell' \text{endif}^{\ell t}} \langle \text{fin}[\text{if } \ell \text{b then } c_1 \text{ else } c_2 \text{ endif}^{\ell t}], \rho \rangle \mid \rho \in \mathcal{E}\}$
$T[c_1; c_2] = T[c_1] \cup T[c_2]$
$T[\text{while } \ell \text{b do } c \text{ end}^{\ell t}] = \{\langle \ell, \rho \rangle \xrightarrow{\ell \text{not b}} \langle \ell', \rho \rangle \mid \rho \in \mathcal{E} \wedge \text{false} = B[\text{b}]\rho\} \cup$ $\{\langle \ell, \rho \rangle \xrightarrow{\ell \text{b}} \langle \text{in}[c], \rho \rangle \mid \rho \in \mathcal{E} \wedge \text{true} = B[\text{b}]\rho\} \cup T[c] \cup$ $\{\langle \ell', \rho \rangle \xrightarrow{\ell' \text{done}^{\ell t}} \langle \text{fin}[\text{while } \ell \text{b do } c \text{ end}^{\ell t}], \rho \rangle \mid \rho \in \mathcal{E}\}$

A *truth-assignment* is a function  $\bar{Y} : \bar{V} \rightarrow \{\text{T}, \text{F}\}$  that assigns to each propositional variable the value true (T) or false (F). Given a formula  $f \in \Omega(I)$ ,  $\bar{Y} \models f$  means that  $\bar{Y}$  satisfies  $f$ , and  $f_1 \models f_2$  is a shorthand for “ $\bar{Y} \models f_1$  implies  $\bar{Y} \models f_2$ ”.  $\Omega(I)$  is ordered by  $f_1 \preceq f_2 \Leftrightarrow f_1 \models f_2$ . Two formulae  $f_1$  and  $f_2$  are logically equivalent, denoted  $f_1 \equiv f_2$  iff  $f_1 \preceq f_2$  and  $f_2 \preceq f_1$ .

The *unit assignment*  $u$  is defined by  $u(x) = \text{T}$  for all  $\bar{x} \in \bar{V}$ . We define the set of positive formulae by  $\text{Pos} = \{f \in \Omega(I) \mid u \models f\}$ . Some obvious examples are  $\text{T}, x_1 \in \text{Pos}$  and  $\text{F}, \neg x_1 \notin \text{Pos}$ .

We can consider the propositional formula  $\psi$  as a conjunction of subformulae  $(\zeta_0 \wedge \dots \wedge \zeta_n)$ . We denote the set of subformulae of  $\psi$  as  $\text{Sub}_\psi$ . Let  $\nabla$  be the least upper bound operator on propositional formula defined by  $\nabla\{\psi_0, \dots, \psi_n\} = \bigwedge\{\text{Sub}_{\psi_0}, \dots, \text{Sub}_{\psi_n}\}$ .  $(\text{Pos}, \preceq, \nabla)$  forms a join semi lattice. Moreover, let  $\ominus : \text{Pos} \times \text{Pos} \rightarrow \text{Pos}$  be a binary operator defined as “simplification” between two propositional formulae:  $\psi_0 \ominus \psi_1 = \bigwedge(\text{Sub}_{\psi_0} \setminus \text{Sub}_{\psi_1})$ . This “simplification” permits us to obtain all the implication in  $\psi_0$  which are not contained in  $\psi_1$ .

## 2.4 Abstract Semantics

Our approach is based on the abstract domain of logic formulae representing dependency between variables (which tracks the propagation of sensitive/insensitive information). The detection of possible information leakages is performed by evaluating formulae on truth-assignment functions. In particular, the analysis involves the following steps:

- Constructs at each program point the propositional formula  $\psi$  through a fixpoint algorithm which represents an over-approximation of variable’s dependencies up to that program point.
- Partitions the variables into public and private privacy levels. Apply a truth-assignment function  $\bar{Y}$  that assigns to each propositional variable the value T (true) or the value F (false) if the corresponding variable is private or

public, respectively. If  $\bar{Y}$  does not satisfy  $\psi$  at all program points, then there could be some information leakages.

The logic formulae, obtained from program's instructions, are in the form:

$$\bigwedge_{0 \leq i \leq n \ 0 \leq j \leq m} \{\bar{x}_i \rightarrow \bar{y}_j\}$$

which means that the values of variable  $\bar{y}_j$  could depend on the values of variable  $\bar{x}_i$ . For instance, the formula  $\bar{y} \rightarrow \bar{x}$  represents variable dependency in assignment statement  $x := y$ ; similarly, in case of conditional statement  $\text{if}(x == 0) \text{ then } y := z$  we obtain the formula  $(\bar{x} \rightarrow \bar{y}) \wedge (\bar{z} \rightarrow \bar{y})$ . Notice that the propositional variable  $\bar{v}$  corresponds to the program variable  $v$ .

Formally, an abstract state  $\sigma^\# \in \Sigma^\# \stackrel{\text{def}}{=} L \times \text{Pos}$  is a pair  $\langle \ell, \phi \rangle$  where  $\phi \in \text{Pos}$  represents the dependencies occurred among program variables up to label  $\ell \in L$ . Given a pair  $\sigma^\# = \langle \ell, \phi \rangle$ , we define  $l(\sigma^\#) = \ell$  and  $r(\sigma^\#) = \phi$ . Let  $BV(c)$ , defined in Table 9, be the set of bound variables in command  $c$ .

The *abstract semantics* of a command  $c$  is defined by  $\bar{T}[c]$ . Similar to the concrete domain, we denote the transition from  $\sigma_1^\#$  to  $\sigma_2^\#$  by  $\sigma_1^\# \rightarrow \sigma_2^\#$ . The abstract semantics in the domain of propositional formulae is defined in Table 10.

Consider two sets of abstract states  $S_1$  and  $S_2$  such that  $S_1 = \{\langle \ell_0^1, \psi_0^1 \rangle, \dots, \langle \ell_n^1, \psi_n^1 \rangle\}$  and  $S_2 = \{\langle \ell_0^2, \psi_0^2 \rangle, \dots, \langle \ell_m^2, \psi_m^2 \rangle\}$ . The partial ordering is defined by  $S_1 \sqsubseteq^\# S_2 \Leftrightarrow n \leq m \wedge \forall i \in [0, n], \ell_i^1 = \ell_i^2 \wedge \forall i \in [0, n], \psi_i^1 \sqsubseteq \psi_i^2$ . Let  $S_0, \dots, S_n \in \wp(\Sigma^\#)$  be sets of abstract states.  $(\wp(\Sigma^\#), \sqsubseteq^\#)$  forms a poset since it is reflexive, antisymmetric, and transitive by basic properties of logic implication. The join operator  $\sqcup^\#$  is defined by:

$$\begin{aligned} \sqcup^\# \{S_0, \dots, S_n\} &= \bigcup (S_0, \dots, S_n) \\ &\cup \{\langle \ell, \psi \rangle \mid \psi = \nabla \{\psi' \mid \langle \ell, \psi' \rangle \in \bigcup (S_0, \dots, S_n)\}\} \\ &\setminus \{\langle \ell, \psi \rangle \in \bigcup (S_0, \dots, S_n) \mid \exists \langle \ell, \psi' \rangle \in \bigcup (S_0, \dots, S_n) \wedge \psi \neq \psi'\} \end{aligned}$$

and the meet operator  $\sqcap^\#$  by:

$$\begin{aligned} \sqcap^\# \{S_0, \dots, S_n\} &= \{\langle \ell, \psi \rangle \in S' \mid S' \in \{S_0, \dots, S_n\} \wedge \\ &\forall i \in [0, n]. \exists \langle \ell, \psi'_i \rangle \in S_i \wedge \psi \sqsubseteq \psi'_i\} \end{aligned}$$

Table 9: *BV* function

$BV(\ell \text{ skip}) = \{\emptyset\}$
$BV(\ell v := \text{exp}) = \{\bar{v}\}$
$BV(c_0; c_1) = BV(c_0) \cup BV(c_1)$
$BV(\text{if } \ell b \text{ then } c_0 \text{ else } c_1 \ell' \text{ endif}) = BV(c_0) \cup BV(c_1)$
$BV(\text{while } \ell b \text{ do } c \ell' \text{ done}) = BV(c)$



Table 10: Abstract semantics

$\begin{aligned} \bar{T}[\text{skip}^{\ell}] &= \{\langle \ell, \psi \rangle \rightarrow \langle \text{fin}[\text{skip}^{\ell}], \psi \rangle\} \\ \bar{T}[\text{v} := \text{exp}^{\ell}] &= \{\langle \ell, \psi \rangle \rightarrow \langle \text{fin}[\text{v} := \text{exp}^{\ell}], \psi_0 \rangle\} \\ \bar{T}[c_0; c_1] &= \bar{T}[c_0] \cup \bar{T}[c_1] \\ \bar{T}[\text{if } \ell \text{b then } c_0 \text{ else } c_1 \text{ }^{\ell} \text{endif}^{\ell}] &= \bar{T}[c_0] \cup \bar{T}[c_1] \cup \\ &\quad \{\langle \ell, \psi \rangle \rightarrow \langle \text{in}[c_0], \psi \rangle\} \cup \{\langle \ell, \psi \rangle \rightarrow \langle \text{in}[c_1], \psi \rangle\} \cup \\ &\quad \{\langle \ell', \psi \rangle \rightarrow \langle \text{fin}[\text{if } \ell \text{b then } c_0 \text{ else } c_1 \text{ }^{\ell} \text{endif}^{\ell}], \psi_1 \rangle\} \\ &\quad \{\langle \ell', \psi \rangle \rightarrow \langle \text{fin}[\text{if } \ell \text{b then } c_0 \text{ else } c_1 \text{ }^{\ell} \text{endif}^{\ell}], \psi_2 \rangle\} \\ \bar{T}[\text{while } \ell \text{b do } c \text{ }^{\ell} \text{done}^{\ell}] &= \bar{T}[c] \cup \{\langle \ell, \psi \rangle \rightarrow \langle \text{in}[c], \psi \rangle\} \cup \\ &\quad \{\langle \ell', \psi \rangle \rightarrow \langle \text{fin}[\text{while } \ell \text{b do } c \text{ }^{\ell} \text{done}^{\ell}], \psi_3 \rangle\} \end{aligned}$
<p>where</p> $\begin{aligned} \psi_0 &= \bigwedge \{\bar{y} \rightarrow \bar{x} \mid \bar{y} \in \bar{V}(\text{exp}) \wedge \bar{y} \neq \bar{x}\} \\ &\quad \bigwedge \{\bar{z} \rightarrow \bar{w} \mid \bar{z} \rightarrow \bar{x}, \bar{x} \rightarrow \bar{w} \in \psi\} \wedge (\psi \ominus \bigwedge \{\bar{y} \rightarrow \bar{x} \mid \bar{y} \in \bar{V} \wedge \bar{x} \notin \bar{V}[\text{exp}]\}) \\ \psi_1 &= \bigwedge \{\bar{y} \rightarrow \bar{x} \mid \bar{y} \in \bar{V}(\text{b}) \wedge \bar{x} \in BV(c_0) \wedge \bar{y} \neq \bar{x}\} \wedge \psi \\ \psi_2 &= \bigwedge \{\bar{y} \rightarrow \bar{x} \mid \bar{y} \in \bar{V}(\text{b}) \wedge \bar{x} \in BV(c_1) \wedge \bar{y} \neq \bar{x}\} \wedge \psi \\ \psi_3 &= \bigwedge \{\bar{y} \rightarrow \bar{x} \mid \bar{y} \in \bar{V}(\text{b}) \wedge \bar{x} \in BV(c) \wedge \bar{y} \neq \bar{x}\} \wedge \psi \end{aligned}$

Basically, the join operator consists in the union of all elements. When two elements have the same label but different formula, the join operator takes the biggest one. Instead, the meet operator considers only the abstract states, with the same label, which are in all elements. In case of different formulae, the meet operator takes the smallest one. By definition join and meet operator are defined for every subset of elements of our domain. Therefore, we can conclude that  $\langle \emptyset(\Sigma^\sharp), \sqsubseteq^\sharp, \emptyset, \Sigma^\sharp, \sqcup^\sharp, \sqcap^\sharp \rangle$  forms a complete lattice.

Let  $I^\sharp[\mathcal{P}] = \{\langle \text{in}[\mathcal{P}], \text{T} \rangle\}$  be the set of possible initial abstract state of program  $\mathcal{P}$ . We define the *abstract semantics* as the set of all finite sets of abstract states, denoted by  $\Sigma^{\star\sharp}$ , reachable during one or more executions, in a finite time. For each element  $S \in \Sigma^{\star\sharp}$  we can denote by  $S^\dagger$  the set of terminal states, defined as  $S^\dagger = \{\sigma_0^\sharp \mid \nexists \sigma_1^\sharp \in S. \sigma_0^\sharp \rightarrow \sigma_1^\sharp \in \bar{T}\}$  and by  $\ell(S)$  all labels of  $S$ . Let  $S_{\sigma_0^\sharp, \sigma_n^\sharp}$  denote a set of states, called *abstract sequence*, that contains a starting state  $\sigma_0^\sharp$  and an ending state  $\sigma_n^\sharp$  such that contains one or more traces from  $\sigma_0^\sharp$  to  $\sigma_n^\sharp$ . We have that  $S_{\sigma_0^\sharp, \sigma_n^\sharp}^\dagger = \{\sigma_n^\sharp\}$ .

We express the abstract semantics in a fixpoint form.

$$\begin{aligned}
\Sigma^{\star\sharp} &= \text{lfp}^\sqsubseteq F^\sharp \text{ where } F^\sharp \in \Sigma^{\star\sharp} \rightarrow \Sigma^{\star\sharp} \\
F^\sharp(X) &\stackrel{\text{def}}{=} \{\sigma^\sharp \mid \sigma^\sharp \in I^\sharp\} \cup \{S_{\sigma_0^\sharp, \sigma_n^\sharp} \mid n \geq 1 \wedge \sigma_0^\sharp \in I^\sharp \wedge S_{\sigma_0^\sharp, \sigma_{n-1}^\sharp} \in X \wedge \sigma_{n-1}^\sharp \rightarrow \sigma_n^\sharp \in \bar{T}\} \cup \\
&\quad \{\sqcup^\sharp \{S_{\sigma_0^\sharp, \sigma_n^\sharp} \mid S_{\sigma_0^\sharp, \sigma_n^\sharp} \in X\}\}
\end{aligned}$$

Table 11: Results of the analysis by Pos domain

Label	Propositional formula
4	$\bar{x} \rightarrow \bar{y}$
5	$\bar{p} \rightarrow \overline{\text{sum}}$
8	$(\bar{x} \rightarrow \bar{y}) \wedge (\bar{p} \rightarrow \overline{\text{sum}}) \wedge (\bar{y} \rightarrow \overline{\text{sum}})$
10	$(\bar{x} \rightarrow \bar{y}) \wedge (\bar{p} \rightarrow \overline{\text{sum}}) \wedge (\bar{x} \rightarrow \overline{\text{sum}})$
12	$(\bar{x} \rightarrow \bar{y}) \wedge (\bar{p} \rightarrow \overline{\text{sum}}) \wedge (\bar{x} \rightarrow \overline{\text{sum}}) \wedge (\bar{y} \rightarrow \overline{\text{sum}})$
14	$(\bar{x} \rightarrow \bar{y}) \wedge (\bar{p} \rightarrow \overline{\text{sum}}) \wedge (\bar{x} \rightarrow \overline{\text{sum}}) \wedge (\bar{y} \rightarrow \overline{\text{sum}}) \wedge (\bar{n} \rightarrow \overline{\text{sum}}) \wedge (\bar{i} \rightarrow \overline{\text{sum}}) \wedge (\bar{i} \rightarrow \bar{n}) \wedge (\bar{k} \rightarrow \overline{\text{sum}}) \wedge (\bar{k} \rightarrow \bar{n})$

*Example 1.* In order to better understand how our dependency analysis works, consider the code in Figure 1 and the program points 4, 5, 8, 10, 12 and 14. When we apply the steps defined above we obtain the propositional formulae in Table 11.

Through our analysis we tracked all the relation between variables. Suppose that variables  $\{\bar{x}, \bar{p}\}$  are private, while all other variables are public. Formally, the correspondent truth-assignment function is defined by  $\bar{Y} = \{\bar{x}, \bar{p} \mapsto \text{T}\} \cup \{\bar{v} \mapsto \text{F} : \bar{v} \in \text{V} \setminus \{\bar{x}, \bar{p}\}\}$ .  $\bar{Y}$  does not satisfy the propositional formulae since in all considered program points there are some public variables that depends on one or more private variables.

Notice that we detect several spurious relations, too. For instance, in contrast with the obtained result, the variable  $\overline{\text{sum}}$  does not depend on  $\bar{n}$ . Indeed at the end of the both branches the variable  $\overline{\text{sum}}$  has always the same value. In Section 3 we will refine the results through the domains combination.

## 2.5 An Instrumented Concrete Domain

To simplify the proof that our concrete and abstract domains from a Galois connection, we introduce another domain, isomorphic to the concrete domain. Let  $\sigma^\diamond \in \Sigma^\diamond = \text{L} \times \text{A}$  be the set of states of this intermediate domain. A pair  $\langle \ell, \mathbf{a} \rangle \in \text{L} \times \text{A}$  represents an action  $\mathbf{a}$  which occurs at program label  $\ell$ . Consider the set  $\Sigma^{\star\diamond}$  which contains all the possible traces of  $\sigma^\diamond$  that can occur during a finite computation. Given  $\Pi_0^\diamond, \Pi_1^\diamond \in \wp(\Sigma^{\star\diamond})$ , we define that  $\Pi_0^\diamond \subseteq \Pi_1^\diamond$  if and only if for each  $\pi_0^\diamond \in \Pi_0^\diamond$  there exists a  $\pi_1^\diamond \in \Pi_1^\diamond$  such that  $\pi_0^\diamond \leq^\diamond \pi_1^\diamond$ . We have that  $\pi_0^\diamond \leq^\diamond \pi_1^\diamond$  if and only if  $\pi_0^\diamond$  is a subsequence of  $\pi_1^\diamond$ . Therefore  $\langle \wp(\Sigma^{\star\diamond}), \subseteq, \emptyset, \Sigma^{\star\diamond}, \cap, \cup \rangle$  forms a lattice. Moreover, we denote by  $\pi^{\diamond+}$  the last state of the sequence.

We can relate  $\wp(\Sigma^\star)$  and  $\wp(\Sigma^{\star\diamond})$  by an abstraction  $\alpha^\diamond \in \wp(\Sigma^\star) \rightarrow \wp(\Sigma^{\star\diamond})$  and a concretization  $\gamma^\diamond \in \wp(\Sigma^{\star\diamond}) \rightarrow \wp(\Sigma^\star)$  function.

Let  $\text{X} = \{\pi_0, \dots, \pi_n\} \in \wp(\Sigma^\star)$  be a set of partial traces and let  $\text{Y} = \{\pi_0^\diamond, \dots, \pi_n^\diamond\} \in \wp(\Sigma^{\star\diamond})$  be a set of sequences of  $\sigma^\diamond$ .

$$\begin{aligned} \alpha^\diamond(\text{X}) &\equiv \{\langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \dots \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \mid \sigma_0 \xrightarrow{\ell_0 \mathbf{a}_0} \dots \xrightarrow{\ell_m \mathbf{a}_m} \sigma_{m+1} \in \text{X}\} \\ \gamma^\diamond(\text{Y}) &\equiv \{\pi \in \wp(\Sigma^\star) \mid \alpha^\diamond(\{\pi\}) \subseteq \text{Y}\} \end{aligned}$$

**Lemma 1.**  $\wp(\Sigma^*) \xleftrightarrow[\alpha^\circ]{\gamma^\circ} \wp(\Sigma^{*\circ})$  forms an isomorphism, that is,  $\gamma^\circ \circ \alpha^\circ = \alpha^\circ \circ \gamma^\circ = id$  (where  $id$  is the identity function).

*Proof.* We have to prove that  $\gamma^\circ \circ \alpha^\circ = \alpha^\circ \circ \gamma^\circ = id$ , where  $id$  is the identity function. Let  $X$  and  $Y$  be elements of  $\wp(\Sigma^{*\circ})$  and  $\wp(\Sigma^*)$  respectively.

$$\begin{aligned} \alpha^\circ(\gamma^\circ(X)) &= \{\langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \dots \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \mid \sigma_0 \xrightarrow{\ell_0 \mathbf{a}_0} \dots \xrightarrow{\ell_m \mathbf{a}_m} \sigma_{m+1} \in \gamma^\circ(X)\} \\ &\quad \text{by definition of } \alpha^\circ \\ &= \{\langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \dots \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \mid \sigma_0 \xrightarrow{\ell_0 \mathbf{a}_0} \dots \xrightarrow{\ell_m \mathbf{a}_m} \sigma_{m+1} \in \{\pi \mid \alpha^\circ(\{\pi\}) \subseteq X\}\} \\ &\quad \text{by definition of } \gamma^\circ \\ &= \{\langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \dots \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \mid \langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \dots \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \in X\} \\ &= X \end{aligned}$$

$$\begin{aligned} \gamma^\circ(\alpha^\circ(Y)) &= \{\pi \in \wp(\Sigma^*) \mid \alpha^\circ(\{\pi\}) \subseteq \alpha^\circ(Y)\} \\ &\quad \text{by definition of } \gamma^\circ \\ &= \{\pi \in \wp(\Sigma^*) \mid \alpha^\circ(\{\pi\}) \subseteq \{\alpha^\circ(\{\pi'\}) \mid \pi' \in Y\}\} \\ &\quad \text{by definition of } \alpha^\circ \\ &= \{\pi \in \wp(\Sigma^*) \mid \pi \in Y\} \\ &= Y \end{aligned}$$

□

Now we define the relation between  $\wp(\Sigma^{*\circ})$  and  $\wp(\Sigma^{*\#})$  by  $\alpha^\#$  and  $\gamma^\#$ .  $\alpha^\# : \wp(\Sigma^{*\circ}) \rightarrow \wp(\Sigma^{*\#})$  is defined by  $\alpha^\#(X) = \sqcup^\# \{\theta(\pi^\circ) \mid \pi^\circ \in X\}$ , where  $\theta : \Sigma^{*\circ} \rightarrow \wp(\Sigma^{*\#})$  is defined as follows.

$$\begin{aligned} \theta(X) &= \{\langle \ell, \psi \rangle \mid \forall \pi \in X. \forall \pi' = \langle \ell_0, \mathbf{a}_0 \rangle \rightarrow \langle \ell_m, \mathbf{a}_m \rangle \leq^\circ \pi : \\ &\quad m \geq 0 \wedge \ell = \ell_m \wedge \psi = \mathbf{f}_0 \wedge \dots \wedge \mathbf{f}_m\} \end{aligned}$$

such that:

1.  $(\forall \langle \ell, \mathbf{v} := \text{exp} \rangle \in \pi' : \forall \langle \ell', \mathbf{v} := \text{exp}' \rangle \in \pi'. \ell' \leq \ell). \exists \mathbf{f}_i = \bar{\mathbf{y}} \rightarrow \bar{\mathbf{v}} : \bar{\mathbf{y}} \in \bar{\mathbf{V}}(\text{exp})$
2.  $\forall (\langle \ell_i, \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{endif} \rangle) \vee (\langle \ell_i, \text{not } \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{endif} \rangle) \leq^\circ \pi^\circ$  which represents an if statement and  $\forall \langle \ell_k, \mathbf{v} := \text{exp}_k \rangle : i < k < j$  exists  $\mathbf{f}_i = \bar{\mathbf{y}} \rightarrow \bar{\mathbf{v}}$  such that  $\bar{\mathbf{y}} \in \bar{\mathbf{V}}(\mathbf{b})$ .
3.  $\forall (\langle \ell_i, \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{done} \rangle) \vee (\langle \ell_i, \text{not } \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{done} \rangle) \leq^\circ \pi^\circ$  which represents a while statement and  $\forall \langle \ell_k, \mathbf{v} := \text{exp}_k \rangle : i < k < j$  exists  $\mathbf{f}_i = \bar{\mathbf{y}} \rightarrow \bar{\mathbf{v}}$  such that  $\bar{\mathbf{y}} \in \bar{\mathbf{V}}(\mathbf{b})$ .

Intuitively, the function  $\theta$  transforms each action (or sequence of actions) in one or more propositional formulae. The easiest (case 1) applies when the action is an assignment statement ( $\mathbf{v} := \text{exp}$ ): we simply obtain the corresponding formula as defined in the transition semantics  $T$ . Instead, for if statements (case 2), we track all the assignment actions that are between if and endif. while statements are treated in a similar way (case 3).

Notice that  $\langle \ell_i, \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{endif} \rangle$  (or  $\langle \ell_i, \text{not } \mathbf{b} \rangle \rightarrow \dots \rightarrow \langle \ell_j, \text{endif} \rangle$ ) represents an if statement if and only if  $\forall (\langle \ell_p, \mathbf{b} \rangle \vee \langle \ell_p, \text{not } \mathbf{b} \rangle) : i < p < j. \exists (\langle \ell_q, \text{endif} \rangle \vee \langle \ell_q, \text{done} \rangle) : p < q < j$  and  $\forall (\langle \ell_q, \text{endif} \rangle \vee \langle \ell_q, \text{done} \rangle) : i < q < j. \exists (\langle \ell_p, \mathbf{b} \rangle \vee \langle \ell_p, \text{not } \mathbf{b} \rangle) : i < p < q$ . Similarly for while statement.

Informally, the pair if and endif (or while and done) is an if (while) statement if and only if between these two actions, there are only assignments or other pairs if-endif or while-done which correspond to nested if and while statements. To better understand, consider the sequence  $\dots \langle \ell_0, \mathbf{b}_0 \rangle \rightarrow \langle \ell_1, \mathbf{b}_1 \rangle \rightarrow \langle \ell_2, \mathbf{v} := \text{exp} \rangle \rightarrow \langle \ell_3, \text{endif} \rangle \dots$ : the pairs  $\langle \ell_0, \mathbf{b}_0 \rangle$  and  $\langle \ell_3, \text{endif} \rangle$  are not an if statement because between these two actions there is  $\langle \ell_1, \mathbf{b}_1 \rangle$ , which does not represent an assignment action neither an if statement.

The concretization function  $\gamma^\# : \wp(\Sigma^{\star\#}) \rightarrow \wp(\Sigma^{\star\circ})$  is defined by  $\gamma^\#(\mathbf{Y}) = \{\pi^\circ \in \Sigma^{\star\circ} \mid \theta(\pi^\circ) \sqsubseteq^\# \mathbf{Y} \wedge l(\pi^\circ) \in \ell(\mathbf{Y}^\dagger)\}$  where  $\mathbf{Y} \in \wp(\Sigma^{\star\#})$ .

**Lemma 2.**  $\theta : \Sigma^{\star\circ} \rightarrow \wp(\Sigma^{\star\#})$  is monotonic:  $x \leq^\circ y \Rightarrow \theta(x) \sqsubseteq^\# \theta(y)$

*Proof.* Let  $x_0 = \{\sigma_0 \rightarrow \dots \rightarrow \sigma_n\}$  and  $x_1 = \{\sigma'_0 \rightarrow \dots \rightarrow \sigma'_m\}$  be two elements of  $\Sigma^{\star\circ}$  such that  $x_0 \leq^\circ x_1$  and consider  $\theta(x_0) = \{\sigma^{\#}_0, \dots, \sigma^{\#}_n\}$  and  $\theta(x_1) = \{\sigma^{\#}_0, \dots, \sigma^{\#}_m\}$ . By the definition of " $\leq^\circ$ " we know that  $n \leq m, \forall i \in [0, n]. \sigma_i = \sigma'_i$ . Therefore, by the definition of  $\theta$ , we have that  $\forall i \in [0, n]. \sigma^{\#}_i = \sigma^{\#}_i$ . Then, by definition of " $\sqsubseteq^\#$ ",  $\theta(x_0) \sqsubseteq^\# \theta(x_1)$ .  $\square$

**Lemma 3.**  $\alpha^\# : \wp(\Sigma^{\star\circ}) \rightarrow \wp(\Sigma^{\star\#})$  is monotonic:  $X \subseteq Y \Rightarrow \alpha^\#(X) \sqsubseteq^\# \alpha^\#(Y)$

*Proof.* Consider  $X_0, X_1 \in \wp(\Sigma^{\star\circ})$  such that  $X_0 \subseteq X_1, \alpha^\#(X_0) = \sqcup^\# \{\theta(\pi^\circ) \mid \pi^\circ \in X_0\}$  and  $\alpha^\#(X_1) = \sqcup^\# \{\theta(\pi^\circ) \mid \pi^\circ \in X_1\}$ . By definition of " $\subseteq$ ",  $\forall \pi^\circ \in X_0, \exists \pi^\circ \in X_1$ . By Lemma 2,  $\theta(\pi^\circ_0) \sqsubseteq^\# \theta(\pi^\circ_1)$  for all  $\pi^\circ_0 \in X_0$  and  $\pi^\circ_1 \in X_1$ . Then we have  $\alpha^\#(X_0) \sqsubseteq^\# \alpha^\#(X_1)$ :  $\alpha^\#(X_1)$  contains all the elements in  $\alpha^\#(X_0)$ .  $\square$

**Lemma 4.**  $\gamma^\# : \wp(\Sigma^{\star\#}) \rightarrow \wp(\Sigma^{\star\circ})$  is monotonic:  $X \sqsubseteq^\# Y \Rightarrow \gamma^\#(X) \subseteq \gamma^\#(Y)$

*Proof.* Consider  $X_0, X_1 \in \wp(\Sigma^{\star\#})$  such that  $X_0 \sqsubseteq^\# X_1, \gamma^\#(X_0) = \{\pi^\circ \in \wp(\Sigma^{\star\circ}) \mid \theta(\pi^\circ) \sqsubseteq^\# X_0 \wedge l(\pi^\circ) \in \ell(X_0^\dagger)\}$  and  $\gamma^\#(X_1) = \{\pi^\circ \in \wp(\Sigma^{\star\circ}) \mid \theta(\pi^\circ) \sqsubseteq^\# X_1 \wedge l(\pi^\circ) \in \ell(X_1^\dagger)\}$ . By definition of " $\sqsubseteq^\#$ " and by Lemma 2, for all  $\pi^\circ_0 \in \gamma^\#(X_0)$  exists  $\pi^\circ_1 \in \gamma^\#(X_1)$ . Therefore  $\gamma^\#(X_0) \subseteq \gamma^\#(X_1)$ .  $\square$

**Lemma 5.**  $\alpha^\# \circ \gamma^\#$  is the identity:  $\alpha^\#(\gamma^\#(X)) = X$

*Proof.* Let  $X$  be an element of  $\wp(\Sigma^{\star\#})$ . By definition of  $\alpha^\#, \alpha^\#(\gamma^\#(X)) = \sqcup^\# \{\theta(\pi^\circ) \mid \pi^\circ \in \gamma^\#(X)\}$ . By definition of  $\gamma^\#, \alpha^\#(\gamma^\#(X)) = \sqcup^\# \{\theta(\pi^\circ) \mid \theta(\pi^\circ) \sqsubseteq^\# X \wedge l(\pi^\circ) \in \ell(X^\dagger)\}$ . Then,  $\alpha^\#(\gamma^\#(X))$  contains the least upper bound of all the abstract traces that have the same last label of  $X$  and that are less or equal than  $X$ . Therefore  $\alpha^\#(\gamma^\#(X)) = X$ .  $\square$

**Lemma 6.**  $\gamma^\# \circ \alpha^\#$  is extensive:  $X \sqsubseteq^\# \gamma^\#(\alpha^\#(X))$

*Proof.* Consider  $X \in \wp(\Sigma^{*\circ})$ . By definition of  $\gamma^\sharp$ ,  $\gamma^\sharp(\alpha^\sharp(X)) = \{\pi^\circ \in \wp(\Sigma^{*\circ}) \mid \theta(\pi^\circ) \sqsubseteq^\sharp \alpha^\sharp(X) \wedge l(\pi^\circ) \in \ell(\alpha^\sharp(X)^+)\}$ . By definition of  $\alpha^\sharp$ ,  $\gamma^\sharp(\alpha^\sharp(X)) = \{\pi^\circ \in \wp(\Sigma^{*\circ}) \mid \theta(\pi^\circ) \sqsubseteq^\sharp \sqcup^\sharp\{\theta(\pi^\circ) \mid \pi^\circ \in X\} \wedge l(\pi^\circ) \in \ell(\alpha^\sharp(X)^+)\}$ . By definition of “ $\sqcup^\sharp$ ”, “ $\sqsubseteq^\sharp$ ” and by Lemma 2,  $X \sqsubseteq^\sharp \gamma^\sharp(\alpha^\sharp(X))$ .  $\square$

**Lemma 7.**  $\wp(\Sigma^{*\circ}) \xleftrightarrow[\alpha^\sharp]{\gamma^\sharp} \wp(\Sigma^{*\sharp})$  is a Galois insertion.

*Proof.*  $\wp(\Sigma^{*\circ})$  and  $\wp(\Sigma^{*\sharp})$  are two complete lattices,  $\gamma^\sharp$  and  $\alpha^\sharp$  are monotonic (Lemma 3 and 4),  $\alpha^\sharp \circ \gamma^\sharp$  is the identity (Lemma 5) and  $\gamma^\sharp \circ \alpha^\sharp$  is extensive (Lemma 6). Therefore  $\wp(\Sigma^{*\circ}) \xleftrightarrow[\alpha^\sharp]{\gamma^\sharp} \wp(\Sigma^{*\sharp})$  is a Galois insertion.  $\square$

Finally, we can express the relation between  $\wp(\Sigma^*)$  and  $\wp(\Sigma^{*\sharp})$  by the composition of above functions,  $\alpha = \alpha^\sharp \circ \alpha^\circ$  and  $\gamma = \gamma^\circ \circ \gamma^\sharp$ . Since the composition of an isomorphism and a Galois insertion is a Galois insertion, we can assert that  $\wp(\Sigma^*) \xleftrightarrow[\gamma^\sharp]{\alpha^\sharp} \wp(\Sigma^{*\sharp})$  is a Galois insertion.

## 2.6 Properties

The aim of information flow analysis is to verify the confidentiality and the integrity of the information in computer programs. An information flow analysis can be carried out by considering different attacker abilities. In this context we consider two different scenarios: when the attacker can read public variables only at the beginning and at the end of the computation, and when the attacker can read public variables after each step of the computation. Note that the attacker, in both cases, knows the source code of the program.

Both the properties and the types of attacker are checked through the definition and the satisfiability of the propositional formulae (Pos) with respect to the truth-assignment function. Let  $\overline{\gamma_\mathcal{P}} : \overline{V} \rightarrow \{T, F\}$  be a truth-assignment function associated with the program  $\mathcal{P}$ . The security properties are modeled by the function definition, while the attacker is modeled by the set of propositional formulae we consider for the satisfiability. For the first case, in which the attacker can read public variables only at the beginning and at the end of the computation, the set of states to consider involves only the terminal states of each sequence ( $\{S \in \Sigma^{*\sharp} \mid \overline{\gamma_\mathcal{P}} \models r(S^+)\}$ ). Whereas in the second case, when the attacker can read public variables at each step of the computation, the set of states to consider involves all the propositional formulae in the sequence ( $\{S \in \Sigma^{*\sharp} \mid \forall \sigma^\sharp \in S : \overline{\gamma_\mathcal{P}} \models r(\sigma^\sharp)\}$ ).

**Confidentiality** Confidentiality refers to limiting information access and disclosure to authorized users. For example, we require when we buy something online that our private data (e.g., credit card number) can be read only by the merchant.

Let  $\Upsilon_{\mathcal{P}} : \mathcal{V} \rightarrow \{\text{L}, \text{H}\}$  be a function which assigns to each variable of program  $\mathcal{P}$  a security class.  $\mathcal{P}$  respects the confidentiality property, if and only if it does not contain any information leakage with respect to the function  $\Upsilon_{\mathcal{P}}$ , i.e., there is no information that moves from private to public variables. To verify this property, we define the corresponding truth-assignment function  $\overline{\Upsilon_{\mathcal{P}}}$  as follows.

$$\overline{\Upsilon_{\mathcal{P}}}(\bar{x}) = \begin{cases} \text{T} & \text{if } \Upsilon_{\mathcal{P}}(x) = \text{H} \\ \text{F} & \text{if } \Upsilon_{\mathcal{P}}(x) = \text{L} \end{cases}$$

**Integrity** By integrity we mean that unauthorized people cannot modify a message.

Let  $\Upsilon_{\mathcal{P}} : \mathcal{V} \rightarrow \{\text{L}, \text{H}\}$  be a function which assigns to each variable of program  $\mathcal{P}$  a security class. The integrity property is verified if and only if public variables do not modify private variables, i.e., there is no information leakage from public variables to private variables. The corresponding truth-assignment function  $\overline{\Upsilon_{\mathcal{P}}}$ , to check this property, is defined as follows.

$$\overline{\Upsilon_{\mathcal{P}}}(\bar{x}) = \begin{cases} \text{T} & \text{if } \Upsilon_{\mathcal{P}}(x) = \text{L} \\ \text{F} & \text{if } \Upsilon_{\mathcal{P}}(x) = \text{H} \end{cases}$$

Notice that it is exactly the opposite of the truth-assignment function for the confidentiality property.

### 3 Combination of Symbolic and Numerical Domains

In this Section, we combine the symbolic propositional formulae domain described above with a numerical domain through reduced product, yielding to a refinement of the results obtained by the dependency analysis. Our modular construction allows to tune efficiency and accuracy changing the numerical domain. For instance, if we use intervals, we will be less precise than by using polyhedra, but we will obtain a more efficient analysis.

Let us briefly recall the main features of some numerical domains already in the literature.

*Intervals* Intervals approximate a set of integers by an interval enclosing all of them. Formally, a set  $V \subseteq \mathbb{Z}$  is approximated with  $[a, b]$  where  $a = \min V$  and  $b = \max V$ . If it is not possible to know precisely the upper and lower bound of a set of integers  $a$  and  $b$  are  $-\infty$  and  $+\infty$ , respectively. This domain is a lattice, and the ordering operator  $\sqsubseteq$  is such that  $[a, b] \sqsubseteq [c, d]$  if and only if the interval  $[a, b]$  is contained by  $[c, d]$ . Therefore the top element is the interval  $[-\infty, +\infty]$  and the bottom element is an interval such that  $a > b$ . This lattice has infinite height and contains infinite ascending chains. So it needs a widening operator. Intervals scale up, but in some cases they are too rough.

*Polyhedra* Convex polyhedra are regions of some  $n$ -dimensional space that are bounded by a finite set of hyperplanes. A convex polyhedron in  $\mathbb{R}^n$  describes a relation between  $n$  quantities. P. Cousot and N. Halbwachs [15] applied the theory of abstract interpretation to the static determination of linear equalities and inequalities among program variables by introducing the use of convex polyhedra as an abstract domain.

We denote by  $\mathbf{v} = (v_0, \dots, v_{n-1}) \in \mathbb{R}^n$  a  $n$ -tuple (vector) of real numbers;  $\mathbf{v} \cdot \mathbf{w}$  denotes the scalar product of vectors  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ ; the vector  $\mathbf{0} \in \mathbb{R}^n$  has all components equal to zero. Let  $\mathbf{x}$  be a  $n$ -tuple of distinct variables. Then  $\beta = (\mathbf{a} \cdot \mathbf{x} \bowtie b)$  denotes a linear constraint, for each vector  $\mathbf{a} \in \mathbb{R}^n$ , where  $\mathbf{a} \neq \mathbf{0}, b \in \mathbb{R}$  and  $\bowtie = \{=, \geq, >\}$ . A linear inequality constraint  $\beta$  defines an affine half-space of  $\mathbb{R}^n$ , denoted by  $con(\{\beta\})$ .

A set  $\mathcal{P} \in \mathbb{R}^n$  is a (convex) polyhedron if and only if  $\mathcal{P}$  can be expressed as the intersection of a finite number of affine half-spaces of  $\mathbb{R}^n$ , i.e., as the solution of a finite set of linear inequality constraints. The set of all polyhedra on the vector space  $\mathbb{R}^n$  is denoted as  $\mathcal{P}_n$ . Let  $\langle \mathcal{P}_n, \subseteq, \emptyset, \mathbb{R}^n, \sqcup, \cap \rangle$  be a lattice of convex polyhedra, where " $\subseteq$ " is the set-inclusion, the empty set and  $\mathbb{R}^n$  as the bottom and top elements, respectively. The binary meet operation returns the greatest polyhedron smaller than or equal to the two arguments, correspond to set intersection, and " $\sqcup$ " is the binary join operation and returns the least polyhedron greater than or equal to the two arguments. This abstract domain has exponential complexity, and it does not scale up in practice.

For more details about polyhedra, many works in literature define abstract domains based on polyhedra as Galois connection [6] and implement this domain [5, 27].

*Octagons* A. Miné introduced Octagons [35] for static analysis by abstract interpretation. The author extended a former numerical domain based on Difference-Bound Matrices [34] and showed practical algorithms to represent and manipulate invariants of the form  $\pm x \pm y \leq c$  (where  $x$  and  $y$  are program variables and  $c$  is a real constant) efficiently. Such invariants describe sets of point that are special kind of polyhedra called *octagons* because they feature at most eight edges in a two dimensional space.

The set of invariants which the analysis discovers is a subset of the ones discovered by Polyhedra, but it is quite efficient. In fact, it infers the invariants with a  $O(n^2)$  worst case memory complexity per abstract state and a  $O(n^3)$  worst case time complexity per abstract operation, where  $n$  is the number of variables in the program.

### 3.1 The Reduced Product

The best way to combine the propositional formulae domain  $\langle \wp(\Sigma^{*\#}), \sqsubseteq^\#, \emptyset, \Sigma^\#, \sqcup^\#, \cap^\# \rangle$  and a numerical domain  $\langle \mathfrak{N}, \sqsubseteq^\mathfrak{N}, \perp^\mathfrak{N}, \top^\mathfrak{N}, \sqcup^\mathfrak{N}, \cap^\mathfrak{N} \rangle$  is by using the reduced product operator [14].

Let  $\wp(\Sigma^*) \xleftarrow[\gamma_0]{\alpha_0} \wp(\Sigma^{*\#})$  and  $\wp(\Sigma^*) \xleftarrow[\gamma_1]{\alpha_1} \mathfrak{N}$  be two Galois connections and let  $\varrho : \wp(\Sigma^{*\#}) \times \mathfrak{N} \rightarrow \wp(\Sigma^{*\#}) \times \mathfrak{N}$  be a reduce operator defined as follows: let  $X \in \wp(\Sigma^{*\#})$

be a set of partial traces, and  $\mathfrak{N} \in \mathfrak{N}$  an element of the numerical domain (a set of intervals, an octagon or a polyhedron). Notice that whatever domain you choose,  $\mathfrak{N}$  can be seen as a set of relations among variables value. The reduce operator  $\varrho$  is defined as  $\varrho(\langle X, \mathfrak{N} \rangle) = \langle X', \mathfrak{N} \rangle$  where

$$\begin{aligned} X' = & \{\sigma_{new}^\# \mid \forall \sigma^\# \in X. l(\sigma_{new}^\#) = l(\sigma^\#) \wedge \\ & \wedge r(\sigma_{new}^\#) = (r(\sigma^\#) \ominus \{\bar{x} \rightarrow \bar{y} \mid y = z \in \mathfrak{N}, z \in \bar{V} \cup \mathbb{Z} \wedge z \neq x\})\} \end{aligned}$$

The reduced operator is aimed at excluding pointless dependencies for all variables which have the same value during the execution, without losing purposeful relations (by the condition “ $x \neq z$ ”). The reduce operator removes from the propositional formulae, contained in  $X$ , the implications which have at the right side a variable that has a constant value. In fact if the variable has a constant value, it cannot depend on other variables.

Then, the reduced product  $D^{\sharp}$  is defined as follows:

$$D^{\sharp} = \{\varrho(\langle X, \mathfrak{N} \rangle) \mid X \in \wp(\Sigma^{*\sharp}), \mathfrak{N} \in \mathfrak{N}\}$$

Consider  $X_0, X_1 \in \wp(\Sigma^{*\sharp}), \mathfrak{N}_0, \mathfrak{N}_1 \in \mathfrak{N}$  and  $\langle X_0, \mathfrak{N}_0 \rangle, \langle X_1, \mathfrak{N}_1 \rangle \in D^{\sharp}$ . Then  $\langle X_0, \mathfrak{N}_0 \rangle \sqsubseteq^{\sharp} \langle X_1, \mathfrak{N}_1 \rangle$  if and only if  $X_0 \sqsubseteq^{\sharp} X_1$  and  $\mathfrak{N}_0 \sqsubseteq^{\mathfrak{N}} \mathfrak{N}_1$ . We define the least upper bound and greatest lower bound operator by  $\langle X_0, \mathfrak{N}_0 \rangle \sqcup^{\sharp} \langle X_1, \mathfrak{N}_1 \rangle = \langle X_0 \sqcup^{\sharp} X_1, \mathfrak{N}_0 \sqcup^{\mathfrak{N}} \mathfrak{N}_1 \rangle$  and  $\langle X_0, \mathfrak{N}_0 \rangle \sqcap^{\sharp} \langle X_1, \mathfrak{N}_1 \rangle = \langle X_0 \sqcap^{\sharp} X_1, \mathfrak{N}_0 \sqcap^{\mathfrak{N}} \mathfrak{N}_1 \rangle$ , respectively.

$\langle D^{\sharp}, \sqsubseteq^{\sharp}, \emptyset, \varrho(\langle \Sigma^{*\sharp}, \mathbb{R}^n \rangle), \sqcup^{\sharp}, \sqcap^{\sharp} \rangle$  forms a complete lattice. In order to better understand the improvements yielded by the combination of the two domains consider the following example.

```
foo () {
  0n = 0; 1x = 1; 2i = 0; 3y = x-1; 4sum = p;
  while (5i <= k) do
    if (6n%2 == 0) then
      7sum = y + p; 8n = n + 1;
    else
      9sum = x + (p-1); 10n = n+3;
    11endif
    12i = i+1;
  13done
} 14
```

Fig. 1: Reduced product example

*Example 2.* Consider the code we introduced in Figure 1. We adopt polyhedra as numerical domain. Below we report the results of two analyses for some program points.



### Polyhedra

$$\begin{array}{l|l}
4 & n = 0; x - 1 = 0; i = 0; y = 0 \\
5 & -p + sum = 0; y = 0; x - 1 = 0; -i + n \geq 0; 3i - n \geq 0; \\
8 & -p + sum = 0; y = 0; x - 1 = 0; -i + n \geq 0; -i + k \geq 0; 3i - n \geq 0; \\
10 & -p + sum = 0; y = 0; x - 1 = 0; -i + n \geq 0; -i + k \geq 0; 3i - n \geq 0; \\
12 & -p + sum = 0; y = 0; x - 1 = 0; -i + n - 1 \geq 0; -i + k \geq 0; \\
& i \geq 0; 3i - n + 3 \geq 0; \\
14 & -p + sum = 0; y = 0; x - 1 = 0; -i + n \geq 0; -i + k - 1 \geq 0; 3i - n \geq 0;
\end{array}$$

### Propositional formula

$$\begin{array}{l|l}
4 & x \rightarrow y \\
5 & p \rightarrow sum \\
8 & (x \rightarrow y) \wedge (p \rightarrow sum) \wedge (y \rightarrow sum) \\
10 & (x \rightarrow y) \wedge (p \rightarrow sum) \wedge (x \rightarrow sum) \\
12 & (x \rightarrow y) \wedge (p \rightarrow sum) \wedge (x \rightarrow sum) \wedge (y \rightarrow sum) \\
14 & (x \rightarrow y) \wedge (p \rightarrow sum) \wedge (x \rightarrow sum) \wedge (y \rightarrow sum) \wedge (n \rightarrow sum) \wedge \\
& (i \rightarrow sum) \wedge (i \rightarrow n) \wedge (k \rightarrow sum) \wedge (k \rightarrow n)
\end{array}$$

When we apply the reduce operator defined above we obtain the following propositional formulas:

$$\begin{array}{l|l}
4 & \text{T} \\
5 & p \rightarrow sum \\
8 & p \rightarrow sum \\
10 & p \rightarrow sum \\
12 & p \rightarrow sum \\
14 & (p \rightarrow sum) \wedge (i \rightarrow n) \wedge (k \rightarrow n)
\end{array}$$

By using the reduce operator we simplified the propositional formulas, removing some implications which could in fact generate false alarms when using the direct product of the domains instead of the reduced product. For instance, in Pos analysis we track the relation  $y \rightarrow sum$ . At the same time, in the numerical analysis, we detect that variable  $sum$  is always equal to  $p$  (namely it is constant). This means that  $y \rightarrow sum$  is a false alarm, hence by the reduce product we may delete it. At the same time, we cannot remove the relation between  $sum$  and  $p$  because it is detected also by the numerical analysis.

## 4 An Extension to Database Query Languages

In this section, we extend the full power of the proposed model to the case of data-intensive applications embedding SQL statements, in order to identify possible leakage of sensitive database information as well. This is particular important as in fact unauthorized leakage often occurs while propagating through database applications accessing and processing them legitimately.

## 4.1 A Motivating Example

Consider the database of Table 12 where customer’s personal information and journey-details are stored in tables “Customer” and “Travel” respectively. On booking a particular flight by a customer, the journey details are added to the table “Travel” and the source-destination distance is added to the corresponding entry in ‘DistanceCovered’ attribute of the table “Customer”. Observe that 10 points on the journey each 100 km are offered which is reflected in the attribute ‘Points’. In addition, a boarding-priority value in the attribute ‘BoardPriority’ is assigned to each journey based on the points acquired by the passenger. This is depicted by procedure `BookFlight()` in program  $\mathcal{P}$  in Figure 2.

Assume that values of the attributes ‘Address’, ‘Age’, ‘DistanceCovered’ and ‘Points’ in table “Customer” are private, whereas the information in Table “Travel” is public. To distinguish from the database attributes, we prefix \$ to the application variables in  $\mathcal{P}$ . Finally, suppose the company has decided to upgrade the customers having more than 50 ‘Points’ to the status of ‘BoardPriority’. This is expressed in  $\mathcal{P}$  by the activation of the `Upgrade()` function.

Table 12: Database  $dB$

(a) Table “Customer”

custID	custName	Address	Age	DistanceCovered	Points
1	Alberto	Athens	56	650	60
2	Matteo	Venice	68	49	0
3	Francesco	Washington	38	972	90
4	Smith	Paris	42	185	10

(b) Table “Travel”

custID	Source	Destination	FlightID	JourneyDate	BoardPriority
1	A	B	F139	26-04-14	2
2	C	D	F28	16-11-13	0
3	A	B	F139	26-04-14	3
4	A	B	F139	26-04-14	1

It is clear from the code that the values of ‘BoardPriority’ in tuples where ‘custID’ are equal to ‘1’ and ‘3’ will be upgraded from 2 to 3 and from 3 and 7 respectively. Therefore, an attacker can easily deduce the exact values of sensitive attribute ‘Points’ in Table “Customer”, by observing the change that occurred in the public attribute ‘BoardPriority’ in Table “Travel”.

The example above clearly shows that sensitive database information may be leaked through database applications when public attribute values depend, directly or indirectly, on private attribute values or private application variable values in the program. For instance, in the given example, the leakage occurs due to the dependence “Points $\rightarrow$  BoardPriority” at program label 19.

---

```

Function BookFlight()
1. $flight=checkAvailability($source, $dest);
2. if($flight ≠ NULL){
3.   $dist=computeDistance($source, $dest);
4.   UPDATE Customer SET DistanceCovered = DistanceCovered +
     $dist WHERE custID=$id;
5.   UPDATE Customer SET Points = Points + 10 ×
     FLOOR($dist/100) WHERE custID=$id;
6.   ResultSet rs = SELECT Points FROM Customer WHERE
     custID=$id;
7.   while(rs.next()){
8.     $point=rs.next().Points;
9.     $priority=getPriority($point);
10.    INSERT INTO Travel(userID, Source, Destination,
      FlightID, JourneyDate, BoardPriority) VALUES
      ($id, $source, $dest, $flight, $date, $priority);}
End of Function BookFlight()

...
...

Function Upgrade()
15. ResultSet rs = SELECT custID, DistanceCovered, Points FROM
    Customer WHERE Points>50;
16. while(rs.next()){
17.   $id=rs.next().custID;
18.   $point=rs.next().Points;
19.   UPDATE Travel SET BoardPriority=BoardPriority +
     ($point-50)/10 WHERE custID=$id;}
End of Function Upgrade()

```

---

Fig. 2: Program  $\mathcal{P}$

## 4.2 Labeled Syntax and Concrete Semantics

The labeled syntax description of the language, depicted in Table 13, includes imperative statements embedding SQL. We express an SQL statement by a tuple  $\langle \text{OP}, \phi \rangle$ , where  $\phi$  is a precondition following first-order logic which is used to identify a set of tuples in the database on which the appropriate operation OP (either select, or insert, or update, or delete) is performed. Each operation represents a set of actions, *e.g.* select operation includes GROUP BY, aggregate functions, ORDER BY, etc. Observe that applications embedding SQL statements involve two distinct sets of variables: application variables  $V_a$  and database variables  $V_d$ . Variables from  $V_d$  appear only in the SQL statements, whereas variables in  $V_a$  may appear in all types of instructions (either SQL or imperative).

We define the action function  $a$  and variable function  $V$  for the language in Tables 14 and 15 respectively.

Lets recall from [23] the notion of *environments* correspond to the variables in  $V_a$  and  $V_d$  respectively.

An *application environment*  $\rho_a \in \mathcal{E}_a$  maps a variable  $x \in \text{dom}(\rho_a) \subseteq V_a$  to its value  $\rho_a(x)$ . So,  $\mathcal{E}_a \triangleq V_a \mapsto \mathcal{D}_{\cup}$  where  $\mathcal{D}_{\cup}$  is the semantic domain for  $V_a$ .

Consider a database as a set of indexed tables  $\{t_i \mid i \in I_x\}$  for a given set of indexes  $I_x$ . A *database environment* is defined by a function  $\rho_d$  whose domain is  $I_x$ , such that for  $i \in I_x$ ,  $\rho_d(i) = t_i$ .

Table 13: Syntax of labeled programs embedding SQL

Constants	
$k \in K$	Set of Constants
$k ::= n \mid s$ where $n \in \mathbb{N}, s \in \text{Strings}$	
Variables	
$v_a \in V_a$	Set of Application Variables
$v_a ::= x \mid y \mid z \mid \dots$	
$v_d \in V_d$	Set of Database Attributes
$v_d ::= a_1 \mid a_2 \mid a_3 \mid \dots$	
Expressions	
$\text{exp} \in E$	Set of Arithmetic Expressions
$\text{exp} ::= k \mid v_d \mid v_a \mid \text{exp}_1 \oplus \text{exp}_2$ where $\oplus \in \{+, -, *, /\}$	
$b \in B$	Set of Boolean Expressions
$b ::= \text{true} \mid \text{false} \mid \text{exp}_1 \odot \text{exp}_2 \mid \neg b \mid b_1 \otimes b_2$ where $\odot \in \{\leq, \geq, ==, >, \neq, \dots\}$ and $\otimes \in \{\vee, \wedge\}$	
SQL Preconditions	
$\tau \in T$	Set of Terms
$\tau ::= k \mid v_d \mid v_a \mid f_n(\tau_1, \tau_2, \dots, \tau_n)$ where $f_n$ is an n-ary function.	
$a_f \in A_f$	Set of Atomic Formulas
$a_f ::= R_n(\tau_1, \tau_2, \dots, \tau_n) \mid \tau_1 == \tau_2$ where $R_n(\tau_1, \tau_2, \dots, \tau_n) \in \{\text{true}, \text{false}\}$	
$\phi \in W$	Set of Pre-conditions
$\phi ::= a_f \mid \neg \phi \mid \phi_1 \otimes \phi_2 \mid \odot v \phi$ where $\otimes \in \{\vee, \wedge\}$ and $\odot \in \{\forall, \exists\}$ and $v \in (V_a \cup V_d)$	
SQL Functions	
$g(\mathbf{exp}) ::= \text{GROUP BY}(\mathbf{exp}) \mid id$ where $\mathbf{exp} = \langle \text{exp}_1, \dots, \text{exp}_n \mid \text{exp}_i \in E \rangle$	
$e ::= \text{DISTINCT} \mid \text{ALL}$	
$s ::= \text{AVG} \mid \text{SUM} \mid \text{MAX} \mid \text{MIN} \mid \text{COUNT}$	
$h(\mathbf{exp}) ::= s \circ e(\mathbf{exp}) \mid \text{DISTINCT}(\mathbf{exp}) \mid id$	
$h(*) ::= \text{COUNT}(*)$ where $*$ represents a list of database attributes denoted by $\mathbf{v}_d$	
$h(\mathbf{u}) ::= \langle h_1(u_1), \dots, h_n(u_n) \rangle$ where $\mathbf{h} = \langle h_1, \dots, h_n \rangle$ and $\mathbf{u} = \langle u_1, \dots, u_n \mid u_i = \text{exp} \vee u_i = * \rangle$	
$f(\mathbf{exp}) ::= \text{ORDER BY ASC}(\mathbf{exp}) \mid \text{ORDER BY DESC}(\mathbf{exp}) \mid id$	
Labeled Commands	
$\ell \in L$	Set of Labels
$q \in Q$	Set of Labeled SQL Statements
$q ::= \text{SELECT} \mid \text{UPDATE} \mid \text{INSERT} \mid \text{DELETE}$	
$\text{SELECT} ::= \langle {}^{\ell_5} \text{assign}(\mathbf{v}_a), {}^{\ell_4} f(\mathbf{exp}'), {}^{\ell_3} e(h(\mathbf{u})), {}^{\ell_2} \phi', {}^{\ell_1} g(\mathbf{exp}), {}^{\ell_0} \phi \rangle$	
$\text{UPDATE} ::= \langle {}^{\ell'} \mathbf{v}_d \stackrel{\text{upd}}{=} \mathbf{exp}, {}^{\ell} \phi \rangle$	
$\text{INSERT} ::= \langle {}^{\ell'} \mathbf{v}_d \stackrel{\text{new}}{=} \mathbf{exp}, {}^{\ell} \text{true} \rangle$	
$\text{DELETE} ::= \langle {}^{\ell'} \text{del}(\mathbf{v}_d), {}^{\ell} \phi \rangle$	
$c \in C$	Set of Labeled Commands
$c ::= {}^{\ell} \text{skip} \mid {}^{\ell} v_a = \text{exp} \mid q \mid c_1; c_2$   if ${}^{\ell} b$ then $c_1$ else $c_2$ ${}^{\ell'}$ endif   while ${}^{\ell} b$ do $c$ ${}^{\ell'}$ done	
$\mathcal{P} ::= c^{\ell}$	Program that ends with label $\ell$ .

Table 14: Definition of Action Function  $a$

---

$a[\text{SELECT}] \stackrel{\text{def}}{=} \{\ell_5 \text{assign}(\mathbf{v}_a), \ell_4 f(\mathbf{exp}'), \ell_3 e(h(\mathbf{u})), \ell_2 \phi', \ell_1 g(\mathbf{exp}), \ell_0 \phi\}$
$a[\text{UPDATE}] \stackrel{\text{def}}{=} \{\ell' \mathbf{v}_d \stackrel{\text{upd}}{=} \mathbf{exp}, \ell \phi\}$
$a[\text{INSERT}] \stackrel{\text{def}}{=} \{\ell' \mathbf{v}_d \stackrel{\text{new}}{=} \mathbf{exp}\}$
$a[\text{DELETE}] \stackrel{\text{def}}{=} \{\ell' \text{del}(\mathbf{v}_d), \ell \phi\}$

---

Given a database environment  $\rho_d$  and a table  $t \in d$ . Assume  $\text{attr}(t) = \{a_1, a_2, \dots, a_k\}$ . So,  $t \subseteq D_1 \times D_2 \times \dots \times D_k$  where  $a_i$  is the attribute corresponding to the typed domain  $D_i$ . A *table environment*  $\rho_t$  for a table  $t$  is defined as a function such that for any attribute  $a_i \in \text{attr}(t)$ ,  $\rho_t(a_i) = \langle \pi_i(l_j) \mid l_j \in t \rangle$ , where  $\pi$  is the projection operator and  $\pi_i(l_j)$  represents  $i^{\text{th}}$  element of the  $l_j$ -th row. In other words,  $\rho_t$  maps  $a_i$  to the ordered set of values over the rows of the table  $t$ .

A *state*  $\sigma \in \Sigma \triangleq \mathbb{L} \times \mathcal{E}_d \times \mathcal{E}_a$  is denoted by a tuple  $\langle \ell, \rho_d, \rho_a \rangle$  where  $\ell \in \mathbb{L}$ ,  $\rho_d \in \mathcal{E}_d$  and  $\rho_a \in \mathcal{E}_a$  are the label of the statement to be executed, the database environment and the application environment respectively.

The set of *states* of a program  $\mathcal{P}$  is, thus, defined as  $\Sigma[\mathcal{P}] \triangleq \mathbb{L}[\mathcal{P}] \times \mathcal{E}_d[\mathcal{P}] \times \mathcal{E}_a[\mathcal{P}]$ , where  $\mathbb{L}[\mathcal{P}]$  is the set of labels in  $\mathcal{P}$ , and  $\mathcal{E}_d[\mathcal{P}]$  and  $\mathcal{E}_a[\mathcal{P}]$  are the sets of database and application environments whose domain is the set of database and application variables in  $\mathcal{P}$  only.

The *labeled transition relation*  $\mathbb{T} : \Sigma \times \mathbf{A} \mapsto \wp(\Sigma)$  specifies which successor states  $\sigma' = \langle \ell', \rho_{d'}, \rho_{a'} \rangle \in \Sigma$  can follow when an action  $\mathbf{a} \in \mathbf{A}$  executes on state  $\sigma = \langle \ell, \rho_d, \rho_a \rangle \in \Sigma$ . We denote a labeled transition by  $\sigma \xrightarrow{\mathbf{a}} \sigma'$  or by  $\langle \ell, \rho_d, \rho_a \rangle \xrightarrow{\mathbf{a}} \langle \ell', \rho_{d'}, \rho_{a'} \rangle$ , or by  $\langle \ell, \rho \rangle \xrightarrow{\mathbf{a}} \langle \ell', \rho' \rangle$  where  $\rho$  and  $\rho'$  represent  $(\rho_d, \rho_a)$  and  $(\rho_{d'}, \rho_{a'})$  respectively.

The *labeled transition semantics*  $\mathbb{T}[\mathcal{P}] \in \wp(\Sigma[\mathcal{P}] \times a[\mathcal{P}] \mapsto \wp(\Sigma[\mathcal{P}]))$  of a program  $\mathcal{P}$  restricts the transition relation to program actions, *i.e.*

$$\mathbb{T}[\mathcal{P}]\sigma = \{\sigma' \mid \sigma \xrightarrow{\mathbf{a}} \sigma' \wedge \mathbf{a} \in a[\mathcal{P}] \wedge \sigma, \sigma' \in \Sigma[\mathcal{P}]\}$$

The *labeled transition semantics* of various commands in database applications can easily be defined from the semantic description reported in [23].

Given a program  $\mathcal{P}$ , let  $l = \{\langle \text{in}[\mathcal{P}], \rho_d, \rho_a \rangle \mid \rho_a \in \mathcal{E}_a \wedge \rho_d \in \mathcal{E}_d\}$  be the set of initial states of  $\mathcal{P}$ . The *partial trace semantics* of  $\mathcal{P}$  can be defined as

$$\mathbb{T}[\mathcal{P}](l) = \text{lf}_{\emptyset}^{\subseteq} F(l) = \bigcup_{i \leq \omega} F^i(l)$$

where  $F(l_0) = \lambda X. l_0 \cup \left\{ \sigma_0 \xrightarrow{\mathbf{a}_0} \dots \xrightarrow{\mathbf{a}_{n-1}} \sigma_n \xrightarrow{\mathbf{a}_n} \sigma_{n+1} \mid \sigma_0 \xrightarrow{\mathbf{a}_0} \dots \xrightarrow{\mathbf{a}_{n-1}} \sigma_n \in X \wedge \sigma_n \xrightarrow{\mathbf{a}_n} \sigma_{n+1} \in \mathbb{T}[\mathcal{P}] \right\}$

Table 15: Definition of Variables Function  $V$ 

$V[\llbracket c \rrbracket]$	$\stackrel{def}{=} \emptyset$
$V[\llbracket v \rrbracket]$	$\stackrel{def}{=} \{v\}$ , where $v \in (V_a \cup V_d)$
$V[\llbracket \mathbf{v} \rrbracket]$	$\stackrel{def}{=} \bigcup_{v_i \in \mathbf{v}} V[\llbracket v_i \rrbracket]$
$V[\llbracket \text{exp}_1 \oplus \text{exp}_2 \rrbracket]$	$\stackrel{def}{=} V[\llbracket \text{exp}_1 \rrbracket] \cup V[\llbracket \text{exp}_2 \rrbracket]$ , where $\oplus \in \{+, -, *, /\}$
$V[\llbracket \mathbf{exp} \rrbracket]$	$\stackrel{def}{=} \bigcup_{\text{exp}_i \in \mathbf{exp}} V[\llbracket \text{exp}_i \rrbracket]$
$V[\llbracket \text{true} \rrbracket]$	$\stackrel{def}{=} \emptyset$
$V[\llbracket \text{false} \rrbracket]$	$\stackrel{def}{=} \emptyset$
$V[\llbracket \text{exp}_1 \oslash \text{exp}_2 \rrbracket]$	$\stackrel{def}{=} V[\llbracket \text{exp}_1 \rrbracket] \cup V[\llbracket \text{exp}_2 \rrbracket]$ , where $\oslash \in \{\leq, \geq, ==, >, \neq, \dots\}$
$V[\llbracket \neg b \rrbracket]$	$\stackrel{def}{=} V[\llbracket b \rrbracket]$
$V[\llbracket b_1 \otimes b_2 \rrbracket]$	$\stackrel{def}{=} V[\llbracket b_1 \rrbracket] \cup V[\llbracket b_2 \rrbracket]$ , where $\otimes \in \{\vee, \wedge\}$
$V[\llbracket f_n(\tau_1, \dots, \tau_n) \rrbracket]$	$\stackrel{def}{=} V[\llbracket \tau_1 \rrbracket] \cup \dots \cup V[\llbracket \tau_n \rrbracket]$ , where $f_n$ is an n-ary function.
$V[\llbracket R_n(\tau_1, \dots, \tau_n) \rrbracket]$	$\stackrel{def}{=} V[\llbracket \tau_1 \rrbracket] \cup \dots \cup V[\llbracket \tau_n \rrbracket]$ , where $R_n(\tau_1, \tau_2, \dots, \tau_n) \in \{\text{true}, \text{false}\}$
$V[\llbracket \tau_1 = \tau_2 \rrbracket]$	$\stackrel{def}{=} V[\llbracket \tau_1 \rrbracket] \cup V[\llbracket \tau_2 \rrbracket]$
$V[\llbracket \neg \phi \rrbracket]$	$\stackrel{def}{=} V[\llbracket \phi \rrbracket]$
$V[\llbracket \phi_1 \otimes \phi_2 \rrbracket]$	$\stackrel{def}{=} V[\llbracket \phi_1 \rrbracket] \cup V[\llbracket \phi_2 \rrbracket]$ , where $\otimes \in \{\vee, \wedge\}$
$V[\llbracket \oslash \phi \rrbracket]$	$\stackrel{def}{=} \{v\} \cup V[\llbracket \phi \rrbracket]$ , where $\oslash \in \{\forall, \exists\}$
$V[\llbracket \text{SELECT} \rrbracket]$	$\stackrel{def}{=} V[\llbracket \mathbf{v}_a \rrbracket] \cup V[\llbracket \mathbf{exp}' \rrbracket] \cup V[\llbracket \mathbf{u} \rrbracket] \cup V[\llbracket \phi' \rrbracket] \cup V[\llbracket \mathbf{exp} \rrbracket] \cup V[\llbracket \phi \rrbracket]$
$V[\llbracket \text{UPDATE} \rrbracket]$	$\stackrel{def}{=} V[\llbracket \mathbf{v}_d \rrbracket] \cup V[\llbracket \mathbf{exp} \rrbracket] \cup V[\llbracket \phi \rrbracket]$
$V[\llbracket \text{INSERT} \rrbracket]$	$\stackrel{def}{=} V[\llbracket \mathbf{v}_d \rrbracket] \cup V[\llbracket \mathbf{exp} \rrbracket]$
$V[\llbracket \text{DELETE} \rrbracket]$	$\stackrel{def}{=} V[\llbracket \mathbf{v}_d \rrbracket] \cup V[\llbracket \phi \rrbracket]$

### 4.3 Abstract Semantics

In case of applications embedding SQL statements, we need to consider two additional dependences, called *database-database* dependence and *program-database* dependence [24]. A *program-database* dependence arises between a database variable and an application variable, where values of the database variable depend on the value of the program variable or vice-versa. A *database-database* dependence arises between two database variables where the values of one depend on the values of the other.

*Example 3.* Consider the database of Table 12 in section 4.1. Consider the following SELECT query:

```
q1 = SELECT Points, AVG(Age) INTO va FROM Customer WHERE Points >=50 GROUP BY
Points HAVING SUM(DistanceCovered)>100 ORDER BY Points
```

Note that we use “INTO v<sub>a</sub>” in q<sub>1</sub> to mention that the result of the query is finally assigned to v<sub>a</sub>, where v<sub>a</sub> is a Record or ResultSet type application variable with fields  $w = \langle w_1, w_2 \rangle$ . The type of  $w_1, w_2$  are same as the return type of ‘Points’, ‘AVG(Age)’ respectively. Recall from Table 13 that the syntax of SELECT statement is defined as:

$$\langle \ell^5 \text{assign}(v_a), \ell^4 f(\mathbf{exp}'), \ell^3 e(h(\mathbf{u})), \ell^2 \phi', \ell^1 g(\mathbf{exp}), \ell^0 \phi \rangle$$

According the syntax defined above,  $q_1$  can be formulated as:

```
q1 = SELECT  $\epsilon(h(u))$  INTO  $v_a(w)$  FROM Customer WHERE  $\phi$  GROUP BY(exp) HAVING
 $\phi'$  ORDER BY ASC(exp')
```

where

- $\phi = Points \geq 50$
- **exp** =  $\langle Points \rangle$
- $g(\mathbf{exp}) = GROUP BY(\langle Points \rangle)$
- $\phi' = (SUM \circ ALL(DistanceCovered)) > 100$
- $h = \langle DISTINCT, AVG \circ ALL \rangle$
- $u = \langle Points, Age \rangle$
- $h(u) = \langle DISTINCT(Points), AVG \circ ALL(Age) \rangle$
- **exp'** =  $\langle Points \rangle$
- $f(\mathbf{exp}') = ORDER BY ASC(\langle Points \rangle)$
- $v_a =$  Record or ResultSet type application variable with fields  $w = \langle w_1, w_2 \rangle$ . The type of  $w_1$  and  $w_2$  are same as the return type of  $DISTINCT(Points)$  and  $AVG \circ ALL(Age)$  respectively.

From  $q_1$  we get the following set of logical formula representing variable dependences:

$$\frac{Points}{DistanceCovered} \rightarrow \overline{v_a.w_1}, \quad \frac{Age}{DistanceCovered} \rightarrow \overline{v_a.w_2}, \quad \frac{Points}{DistanceCovered} \rightarrow \overline{v_a.w_2}$$

Below we depict variable dependences in other SQL commands.

```
q2=UPDATE Customer SET DistanceCovered = $y + 150 WHERE custID=2
/* where $y is an application variable. */
```

The logical formula obtained from  $q_2$  are:  $\overline{custID} \rightarrow \overline{DistanceCovered}, \overline{\$y} \rightarrow \overline{DistanceCovered}$ .

```
q3 = INSERT INTO Travel(custID,Source,Destination,FlightID, JourneyDate,BoardPriority)
VALUES (5,"D","E","F34", $y,$z)
/* where $y and $z are application variables. */
```

The logical formula obtained from  $q_3$  are:  $\overline{\$y} \rightarrow \overline{JourneyDate}, \overline{\$z} \rightarrow \overline{BoardPriority}$ .

```
q4 = DELETE FROM Customer WHERE Age > 60
```

The logical formula obtained from  $q_4$  are:

$$\frac{Age}{Age} \rightarrow \overline{custID}, \frac{Age}{Age} \rightarrow \overline{custName}, \quad \frac{Age}{Age} \rightarrow \overline{Address}$$

$$\frac{Age}{Age} \rightarrow \overline{Age}, \quad \frac{Age}{Age} \rightarrow \overline{DistanceCovered}, \frac{Age}{Age} \rightarrow \overline{Points}$$

The dependences above indicate explicit-flow of information. An example of implicit-flow that may occur in case of our application is, for instance, when manipulation of any public database information is performed under the control statements involving high variables.

Table 16 depicts abstract labeled transition semantics of various statements in database applications. The abstract semantics of the program is obtained by fix-point computation over the abstract domain.

Table 16: Definition of Abstract Transition Function  $\bar{T}$

$\bar{T}[\text{SELECT}]$ $\stackrel{def}{=} \bar{T}[\langle \ell_5 \text{assign}(v_a), \ell_4 f(\mathbf{exp}'), \ell_3 e(h(u)), \ell_2 \phi', \ell_1 g(\mathbf{exp}), \ell_0 \phi \rangle]$ $\stackrel{def}{=} \langle \langle \ell_0, \psi \rangle \xrightarrow{\text{SELECT}} \langle \text{fin}[\text{SELECT}], \psi' \rangle \rangle$ <p style="margin-left: 20px;">where <math>\psi' = \bigwedge \{ \bar{v} \rightarrow \bar{v}_a.\bar{w}_i \mid \bar{v} \in (\bar{V}[\phi] \cup \bar{V}[\mathbf{exp}] \cup \bar{V}[\phi'] \cup \bar{V}[\mathbf{exp}']) \wedge \bar{v}_a.\bar{w}_i \in \bar{v}_a.\bar{w} \wedge \bar{v} \neq \bar{v}_a.\bar{w}_i \} \wedge</math>  <math>\bigwedge \{ \bar{v}_i \rightarrow \bar{v}_a.\bar{w}_i \mid \bar{v}_i \in \bar{V}[u_i] \wedge u_i \in u \wedge \bar{v}_a.\bar{w}_i \in \bar{v}_a.\bar{w} \} \wedge (\psi \ominus \bigwedge \{ \bar{v} \rightarrow \bar{v}_a.\bar{w}_i \mid \bar{v} \in \bar{V} \wedge \bar{v}_a.\bar{w}_i \in \bar{v}_a.\bar{w} \})</math></p>
$\bar{T}[\text{UPDATE}]$ $\stackrel{def}{=} \bar{T}[\langle \ell' \mathbf{v}_d \stackrel{upd}{=} \mathbf{exp}, \ell' \phi \rangle]$ $\stackrel{def}{=} \langle \langle \ell, \psi \rangle \xrightarrow{\text{UPDATE}} \langle \text{fin}[\text{UPDATE}], \psi' \rangle \rangle$ <p style="margin-left: 20px;">where <math>\psi' = \bigwedge \{ \bar{v}_1 \rightarrow \bar{v}_2 \mid \bar{v}_1 \in \bar{V}[\phi] \wedge \bar{v}_2 \in \bar{v}_d \} \wedge</math>  <math>\bigwedge \{ \bar{v}_i \rightarrow \bar{v}_j \mid \bar{v}_i \in \bar{V}[\text{exp}_i] \wedge \text{exp}_i \in \mathbf{exp} \wedge \bar{v}_j \in \bar{v}_d \} \wedge \psi</math></p>
$\bar{T}[\text{INSERT}]$ $\stackrel{def}{=} \bar{T}[\langle \ell' \mathbf{v}_d \stackrel{new}{=} \mathbf{exp}, \ell' \text{true} \rangle]$ $\stackrel{def}{=} \langle \langle \ell, \psi \rangle \xrightarrow{\text{INSERT}} \langle \text{fin}[\text{INSERT}], \psi' \rangle \rangle$ <p style="margin-left: 20px;">where <math>\psi' = \bigwedge \{ \bar{v}_i \rightarrow \bar{v}_j \mid \bar{v}_i \in \bar{V}[\text{exp}_i] \wedge \text{exp}_i \in \mathbf{exp} \wedge \bar{v}_j \in \bar{v}_d \} \wedge \psi</math></p>
$\bar{T}[\text{DELETE}]$ $\stackrel{def}{=} \bar{T}[\langle \ell' \text{del}(\mathbf{v}_d), \ell' \phi \rangle]$ $\stackrel{def}{=} \langle \langle \ell, \psi \rangle \xrightarrow{\text{DELETE}} \langle \text{fin}[\text{DELETE}], \psi' \rangle \rangle$ <p style="margin-left: 20px;">where <math>\psi' = \bigwedge \{ \bar{v}_1 \rightarrow \bar{v}_2 \mid \bar{v}_1 \in \bar{V}[\phi] \wedge \bar{v}_2 \in \bar{v}_d \} \wedge \psi</math></p>

#### 4.4 Enhancing the analysis

The dependences that we considered so far are syntax-based, and may yield false positives in the analysis. For instance, let us consider the database in Table 17 and the query  $q_5$ .

$$q_5 = \text{SELECT } Type \text{ INTO } v_a \text{ FROM } Emp, Job \text{ WHERE } Sal = BASIC + (BASIC * (DA/100)) + (BASIC * (HRA/100))$$

The following logical formulae representing PD-dependences exist in  $q_5$ :

$$\psi_5 = \overline{Sal} \rightarrow \overline{v_a.w_1}, \overline{BASIC} \rightarrow \overline{v_a.w_1}, \overline{DA} \rightarrow \overline{v_a.w_1}, \overline{HRA} \rightarrow \overline{v_a.w_1}, \overline{Type} \rightarrow \overline{v_a.w_1}$$

Assuming 'Sal', 'BASIC', 'HRA', 'DA' are private and at least one employee in each job-type must exist, we see that although syntactic PD-dependences above indicating the presence of information leakage, but in practice nothing about these secrets is leaked through  $v_a.w_1$ .

Here is an another example of PD-dependence that is indicating false alarm on leakage: consider the code  $\{ \$x = 4 * \$w * \log 2; \text{UPDATE } t \text{ SET } a = a + \$x; \}$ . Assuming  $\$x$  is private and  $\$w, a$  are public, we see that the dependence  $\$x \rightarrow \bar{a}$  generates false alarm as because  $\$x$  is always equal to 0.

To remove all such false alarms and to increase the accuracy of the analysis, we analyze programs by using the semantic-based abstract interpretation framework.

Consider an abstract domain  $\mathfrak{N}$  where numerical attributes and numerical application variables are abstracted by the *domain of intervals*<sup>8</sup>. The abstraction

<sup>8</sup> For other type of variables, the abstraction function represents identity function.



Table 17: Database  $dB$ (a) Table “ $Emp$ ”

$ID$	$Name$	$Sal$
1	Alberto	1110
2	Matteo	1638
3	Francesco	2255
4	Smith	1840

(b) Table “ $Job$ ”

$Type$	$Rank$	$BASIC$	$HRA$	$DA$
Security	S1	800	20	65
Security	S2	600	20	65
Security	S3	520	15	65
Technical	T1	1000	25	75
Technical	T2	920	25	75
Technical	T3	880	20	75
Technical	T4	840	20	75
Admin	A1	1240	25	80
Admin	A2	1100	25	80

yields an abstract query  $q_6^\#$  corresponding to  $q_6$  and an abstract database depicted in Table 18.

Table 18: Database  $dB^\#$ (a) Table “ $Emp^\#$ ”

$ID^\#$	$Name^\#$	$Sal^\#$
1	Alberto	[1110, 1110]
2	Matteo	[1638, 1638]
3	Francesco	[2255, 2255]
4	Smith	[1840, 1840]

(b) Table “ $Job^\#$ ”

$Type^\#$	$Rank^\#$	$BASIC^\#$	$HRA^\#$	$DA^\#$
Security	[S1, S3]	[520, 800]	[15, 20]	[65, 65]
Technical	[T1, T4]	[840, 1000]	[20, 25]	[75, 75]
Admin	[A1, A2]	[1100, 1240]	[25, 25]	[80, 80]

$$q_6^\# = \text{SELECT}^\#_{Type^\#} \text{INTO}^\#_{v_a^\#} \text{FROM}^\#_{Emp^\#, Job^\#} \text{WHERE}^\#_{Sal^\# =^\# BASIC^\# + (BASIC^\# * (DA^\# / [100, 100])) + (BASIC^\# * (HRA^\# / [100, 100]))}$$

The right-hand side expression of the condition in  $\text{WHERE}^\#$  is evaluated to abstract values [936, 1480], [1638, 2000], and [2255, 2542] respectively corresponding to the three abstract tuples in “ $Job^\#$ ”. Observe that, according to the assumption that at least one employee must exist in each job-type, there exist at least one ‘ $Sal^\#$ ’ in “ $Emp^\#$ ” for which “ $Sal^\# =^\# [936, 1480]$ ” is true, according to the following:

$$[l_i, h_i] =^\# [l_j, h_j] \triangleq \begin{cases} \text{true} & \text{if } (l_i \geq l_j \wedge h_i \leq h_j) \\ \text{false} & \text{if } h_i < l_j \vee l_i > h_j \\ \top & \text{otherwise} \end{cases}$$

Similar for “ $Sal^\# =^\# [1638, 2000]$ ” and “ $Sal^\# =^\# [2255, 2542]$ ”. Therefore, the evaluation of  $q_6$  on  $dB$  always gives the same result *w.r.t.* the property  $\aleph$ , irrespective of the states of “ $Emp$ ”.

We can perform similar analysis of the code  $\{\$x = 4 * \$w * \log 2; \text{UPDATE } t \text{ SET } a = a + \$x;\}$  in the *domain of intervals*, yielding to “no update” of the values in public attribute  $a$ .

The interaction of the logical and numerical domains can be formalized by using the reduced Product  $D^{\sharp}$  as follows:

$$D^{\sharp} = \{\varrho(\langle X, \mathfrak{R} \rangle) \mid X \in \wp(\Sigma^{*\sharp}), \mathfrak{R} \in \mathfrak{N}\}$$

where  $\varrho(\langle X, \mathfrak{R} \rangle) = \{\langle \ell_i, \psi_k \rangle \mid \langle \ell_i, \psi_j \rangle \in X \wedge \psi_k = (\psi_j \ominus \{\bar{v}_1 \rightarrow \bar{v}_2 \mid y \in \gamma(\mathfrak{R})\})\}$ .

In the example above, by analyzing  $q_6$  in the abstract domain  $\mathfrak{N}$  where numerical variables are abstracted by the *domain of intervals*, we see that the value of  $v_a.w_1$  generated by  $q_6$  is always constant throughout the program execution *w.r.t.*  $\mathfrak{N}$ . As  $\psi_6 \in \text{Pos}$  and  $v_a.w_1 \in \mathfrak{N}$ , the reduced product operator  $\varrho$  removes from  $\psi_6$  all dependences in the form “ $x \rightarrow v_a.w_1$ ” (that are representing false alarms), and makes the analysis more accurate and efficient.

## 5 Implementing the Analysis in Sails

In this section, we present `Sails`. The tool is an instance of the generic analyzer `Sample`. This is why we discuss the main issues we have to solve in order to deal with information leakage analysis within `Sample`.

### 5.1 Sample

`Sample` (Static Analyzer of Multiple Programming Languages) is a generic analyzer based on the abstract interpretation theory. Relying on compositional analyses, `Sample` can be plugged with different heap abstractions, approximations of other semantic information (e.g., numeric domains or information flow), properties of interest, and languages. Several heap analyses, semantic and numerical domains have been already plugged. The analyzer works on an intermediate language called `Simple`. Up to now, `Sample` supports the compilation of `Scala` and `Java` bytecode to `Simple`.

Figure 3 depicts the overall structure of `Sample`. Source code programs are compiled to `Simple`. A fixpoint engine receives a heap analysis, a semantic domain, and a control flow graph (whose blocks are composed by a sequence of `Simple` statements), and it produces an abstract result over the control flow graph of each method. This result is passed to a property checker that produces some output (e.g., warnings) to the user. The integration of an analysis in `Sample` allows one to take advantage of all aspects not strictly related to the analysis but that can improve its final precision (e.g., heap or numerical abstractions). For instance, `Sample` is interfaced with the `Apron` library [28] and contains a heap analysis based of `TVLA` [39].

### 5.2 Heap Abstraction

In `Sample` heap locations are approximated by abstract heap identifiers. While the identifiers of program variables are fixed and represent exactly one concrete variable, the abstract heap identifiers may represent several concrete heap

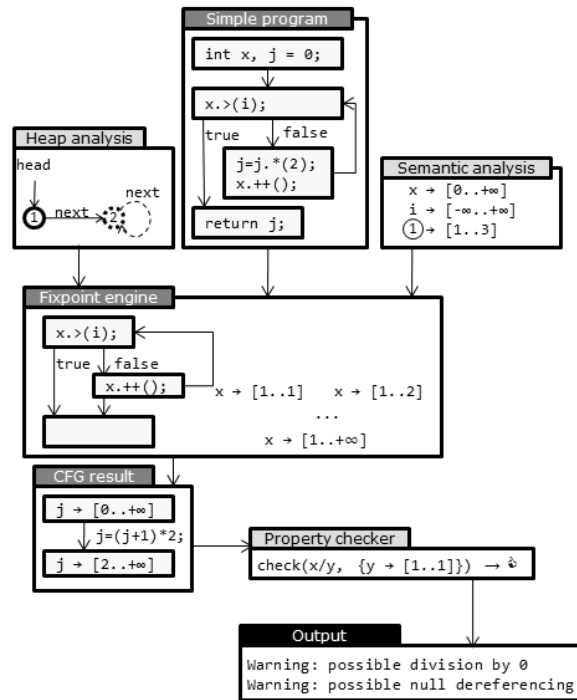


Fig. 3: The structure of Sample

locations (e.g., if they summarize a potentially unbounded list), and they can be merged and split during the analysis. In particular we have to support (i) assignments on summary heap identifiers, and (ii) renaming of identifiers.

In order to preserve the soundness of Sails, we have to perform weak assignments on summary heap identifiers. Since a summary abstract identifier may represent several concrete heap locations and only one of them would be assigned in one particular execution, we have to take the upper bound between the assigned value, and the old one.

The heap abstraction could require to rename, summarize or split existing identifiers. This information is passed through a replacement function  $rep : \wp(\text{Id}) \rightarrow \wp(\text{Id})$ , where  $\text{Id}$  is the set containing all heap identifiers. For instance, in TVLA two abstract nodes represented by identifiers  $a_1$  and  $a_2$  may be merged to a summary node  $a_3$ , or a summary abstract node  $b_1$  may be splitted to  $b_2$  and  $b_3$ . Our heap analysis will pass  $\{a_1, a_2\} \mapsto \{a_3\}$  and  $\{b_1\} \mapsto \{b_2, b_3\}$  to Sails in these cases, respectively. Given a single replacement  $S_1 \mapsto S_2$ , Sails removes all subformulae dealing with some of the variables in  $S_1$ , and for each removed subformula  $s$  it inserts a new subformula  $s'$  renaming each of the variables in  $S_1$

to each of the variables in  $S_2$ . Formally:

$$\begin{aligned}
& \text{rename} : (\text{Pos} \times (\wp(\text{ld}) \rightarrow \wp(\text{ld}))) \rightarrow \text{Pos} \\
& \text{rename}(\sigma^\#, \text{rep}) = \{(i'_1, i'_2) : (i_1, i_2) \in \sigma^\# \wedge \\
& \quad i'_1 = \begin{cases} i_1 & \text{if } \nexists R_1 \in \text{dom}(\text{rep}) : i_1 \in R_1 \\ k_1 & \text{if } \exists R_1 \in \text{dom}(\text{rep}) : i_1 \in R_1 \wedge k_1 \in \text{rep}(R_1) \end{cases} , \\
& \quad i'_2 = \begin{cases} i_2 & \text{if } \nexists R_2 \in \text{dom}(\text{rep}) : i_2 \in R_2 \\ k_2 & \text{if } \exists R_2 \in \text{dom}(\text{rep}) : i_2 \in R_2 \wedge k_2 \in \text{rep}(R_2) \end{cases} \}
\end{aligned}$$

### 5.3 Propositional Formulae

We have to introduce some slight modifications on the domain for information leakage analysis described in Section 2 to work with object oriented languages. We can consider a propositional formula  $\phi$  as a conjunction of subformulae ( $\zeta_0 \wedge \dots \wedge \zeta_n$ ). In the implementation, each subformula is an implication between two identifiers. Then we represent a subformula as a pair of identifiers and a formula as a set of subformulae. Consider the statement  $\text{if}(x > 0) y = z$ . The formula obtained after the analysis of this statement is represented by the set  $\{(\bar{y}, \bar{z}), (\bar{x}, \bar{y})\}$ , where we denote the identifier of the variable  $u$  by  $\bar{u}$ . The order relation “ $\leq$ ” is defined by the subset relation ( $\phi_0 \leq \phi_1 \Leftrightarrow \phi_0 \subseteq \phi_1$ ).

Consequently, in the implementation the set of propositional variables  $\bar{V}$  consists in the set of identifier  $\text{ld}$ , a single propositional formula is represented by  $\wp(\text{ld} \times \text{ld})$  and an abstract state  $\sigma^\# \in \Sigma^\#$  is a conjunction of propositional formulae represented by  $\wp(\wp(\text{ld} \times \text{ld}))$ .

### 5.4 Implicit Flow Detection

An implicit information flow occurs when there is an information leakage from a variable in a condition to a variable assigned inside a block dependent on that condition. For instance, in  $\text{if}(x > 0) y = z$ ; there is an explicit flow from  $z$  to  $y$ , and an implicit flow from  $x$  to  $y$ . To record these relations we relate the variables in the conditions to the variables that have been assigned in the block. When we join two blocks coming from the same condition, we discharge all implicit flows on the abstract state. Observe that **Sails** does not support all cfgs that can be represented in **Sample** but only the ones coming from structured programs, i.e., that corresponds to programs with **if** and **while** statements and not with arbitrary jumps like **goto**.

### 5.5 Property

An information flow analysis can be carried out by considering different attacker abilities. We implemented two scenarios: when the attacker can read public variables only at the beginning and at the end of the computation, and when the attacker can read public variables after each step of the computation<sup>9</sup>.

<sup>9</sup> Notice that, as in [47], we assume that the attacker, in both cases, knows the source code of the program.

Moreover, we implemented two security properties for each attacker: secrecy (i.e., information leakage analysis) and integrity.

The verification of these properties happens after the computation of the analysis and the declaration of private variables (at run time, by a text files writing the variables name or by a graphical user interface selecting the variables in a list).

## 5.6 Numerical Analysis

The information flow analysis is based on the reduced product of a dependency and a numerical analysis. Thanks to the compositional structure of `Sample`, we can plug `Sails` with different numerical domains. In particular, `Sample` supports the `Apron` library. In this way, we can combine `Sails` with all numerical domains contained in `Apron` (namely, `Polka`, the `Parma Polyhedra Library`, `Octagons`, and a deep implementation of `Intervals`).

In addition, we can apply different heap abstractions. For instance, if we are not interested to the heap structure, we can use a less accurate domain that approximates all heap locations with one unique summary node, as we will do in Section 6.2.

## 5.7 Complexity of the Analysis

The complexity of variables dependency analysis showed in Section 2 is strictly correlated to the complexity of propositional formulae. Logical domains, in literature, are widely treated and generally, the logical equivalence of two boolean expression is a co-NP-complete problem. However, this complexity issue may not matter much in practice because the size of the set of variables appearing in the program is reasonably small. Hence, on the one hand, work with propositional formulae requires the solving of a co-NP-complete problem, while on the other hand, in many frameworks (included our system), `Pos` only deal with the variables appearing in the programs, reducing in this way the complexity. Generally, it is possible to increase the efficiency of the computation using the *binary decision diagrams* (BDDs) for the implementation of propositional formulae. For more information about binary decision diagrams see [1].

The simplification adopted in the implementation, i.e. the definition of “ $\preceq$ ” by the subset relation ( $\phi_0 \preceq \phi_1 \Leftrightarrow \phi_0 \subseteq \phi_1$ ), permits to decrease the complexity. In fact, decreasing the precision of the analysis, we can compare two propositional formulae in polynomial time.

About polyhedra analysis, the complexity is well and completely treated in many works [5] and heavily depends on its implementation. For example many implementations, e.g. `Polylib` and `New Polka`, use matrices of coefficients, that cannot grow dynamically, and the worst case space complexity of the methods employed is exponential. In `PPL` library, instead, all data structures are fully dynamic and automatically expanded (in amortized constant time) ensuring the best use of available memory. Comparing the efficiency of polyhedra libraries is not a simple task, because the pay-off depends on the targeted applications: in [5] the authors presented many test results about it.

The complexity of reduced product, and more precisely of reduction operator presented in Section 3.1, is strictly connected with the complexity of the operations on the domains we combine.

## 6 Experimental Results

In this section, we present the experimental results of *Sails*. First of all, we present the results in terms of precision when we analyze a case study involving recursive data structures. Then, we present the results obtained when applying *Sails* to the *SecuriBench-micro* suite.

### 6.1 Case Study

Consider the Java code in Figure 4. Class `ListWorkers` models a list of workers of an enterprise. Each node contains the salary earned by the worker, and some other data (e.g., name and surname of the person). Method `updateSalaries` is defined as well. It receives a list of employees and a list of managers. These two lists are supposed to be disjoint. First method `updateSalaries` computes the maximal salary of an employee. Then it traverses the list of managers updating their salary to the maximal salary of employees if manager's salary is less than that.

Usually managers would not like to leak information about their salary to employees (secrecy property). This property could be expressed in *Sails* specifying that we do not want to have a flow of information *from managers to employees*. More precisely, we want to prove the absence of information leakage from the content of field `salary` of any node reachable from managers to any node reachable from employees.

We combine *Sails* with a heap analysis that approximates all objects created by a program point with a single abstract node [20]. We start the analysis of method `updateSalaries` with an abstract heap in which lists `managers` and `employees` are abstracted with a summary node and they are disjoint. Figure 5 depicts the initial state, where `n2` and `n4` contains the salary values of the `ListWorkers` `n1` and `n3`, respectively. In the graphic representation we adopt dotted circles to represent summary nodes, rectangles to represent local variables, and edges between nodes to represent what is pointed by local variables or fields of objects. Note that the structure of these two lists does not change during the analysis of the program, since method `updateSalaries` does not modify the heap structure.

*Sails* infers that, after the first while loop at line 15, there is a flow of information from `n2` to `maxSalary`. This happens because variable `it` points to `n1` before the loop (because of the assignment at line 9), and it iterates following field `next` (obtaining always the summary node `n1`) perhaps assigning the content of `it.salary` (that is, node `n2`) to `maxSalary`. Therefore, at line 15 we have the propositional formula  $n2 \rightarrow \text{maxSalary}$ .

Then `updateSalaries` traverses the managers list. For each node, it could assign `maxSalary` to `it.salary`. Similarly to what happened in the previous loop,

```

class ListWorkers {
    int salary;
    ListWorkers next;
    ...
}

public void updateSalaries
(ListWorkers employees, ListWorkers managers){
    int maxSalary = 0;
    ListWorkers it=employees;
    while(it!=null) {
        if(it.salary>maxSalary)
            maxSalary=it.salary;
        it=it.next;
    }
    it=managers;
    while(it!=null) {
        if(it.salary < maxSalary)
            it.salary=maxSalary;
        it=it.next;
    }
}

```

Fig. 4: A motivating example

variable `it` points to `n3` before and inside the loop, since field `next` always points to the summary node `n3`. Therefore the assignment at line 18 could potentially affect only node `n4`. For this reason, `Sails` discovers a flow of information from `maxSalary` to `n4`, represented by the propositional formula  $\text{maxSalary} \rightarrow n4$ .

At the end of the analysis, `Sails` soundly computes that  $(n2 \rightarrow \text{maxSalary}) \wedge (\text{maxSalary} \rightarrow n4)$ . By the transitive property, we know that there could be a flow of information from `n2` to `n4`, that is, from `employees` to `managers`. This flow is allowed by our security policy. On the other hand, we also discovered that there is no information leakage from list `managers` to list `employees`, since `Sails` does not contain any propositional formula with this flow. Therefore `Sails` proves that this program is safe.

“Noninterference of programs essentially means that a *variable* of confidential (high) input does not cause a variation of public (low) output”[38]. Thanks to the combination between a heap abstraction and an abstract domain tracking information flow, `Sails` deals directly with the structure of the heap, extending the concept of noninterference from variables to portions of the heap represented by abstract nodes. This opens a new scenario since we can prove that a whole data structure does not interfere with another one, as we have done in this example. As far as we know, `Sails` is the only tool that performs a noninterference analysis over a heap abstraction, and therefore it can prove properties like “there is no

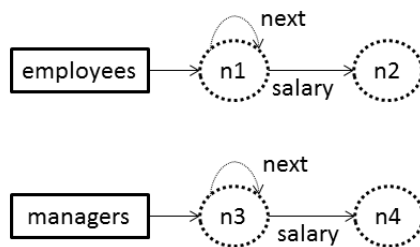


Fig. 5: The initial state of the heap abstraction

information flow from the nodes reachable from  $v_1$  to the nodes reachable from  $v_2$ ".

## 6.2 Benchmarks

A well-established way of studying the precision and the efficiency of information flow analyses is the SecuriBench-micro suite [45]. We applied Sails to this test suite; the description and the results of these benchmarks are reported in Table 19. Column *fa* reports if the analysis did not produce any false alarm. We combined Sails with a really rough heap abstraction that approximates all concrete heap locations with one abstract node. Sails detected all information leakages in all tests, but in three cases (Pred1, Pred6 and Pred7) it produced false alarms. This happens because Sails abstracts away the information produced when testing to true or false boolean conditions in *if* or *while* statements.

Since these benchmarks cover only problems with explicit flows, we performed further experiments using some Jif [36] case studies. The results are reported in Table 20: we discovered all flows without producing any false alarm.

These results allow us to conclude that Sails is precise, since in 90% of the cases (28 out of 31 programs) it does not produce any false alarm.

About the performances, the analysis of all case studies takes 1.092 seconds (0.035 sec per method in average) without combining it with a numerical domain. When we combine it with Intervals it takes 3.015 seconds, whereas it takes 6.130 seconds in combination with Polka. All tests are performed using a MacBook Pro Intel Core 2 Duo 2.53 GHz with 4 GB of RAM memory. Therefore the experimental results underline the efficiency of Sails as well.

## 7 Related Work

In a security-typed language Volpano, Irvine and Smith [46] were the first ones to develop a type system to enforce information flow policies, where a type is inductively associated at compile-time with program statements in such a way that well-typed programs satisfy the non-interference property. The authors



Table 19: SecuriBench-micro suite

Name	Description	fa
Aliasing1	Simple aliasing	✓
Aliasing2	Aliasing false positive	✓
Basic1	Very simple XSS	✓
Basic2	XSS combined with a conditional	✓
Basic3	Simple derived integer test	✓
Basic5	Test of derived integer	✓
Basic6	Complex test of derived integer	✓
Basic8	Test of complex conditionals	✓
Basic9	Chains of value assignments	✓
Basic10	Chains of value assignments	✓
Basic11	A simple false positive	✓
Basic12	A simple conditional	✓
Basic18	Protect against simple loop unrolling	✓
Basic28	Complicated control flow	✓
Pred1	Simple if(false) test	✗
Pred2	Simple correlated tests	✓
Pred3	Simple correlated tests	✓
Pred4	Test with an integer variable	✓
Pred5	Test with a complex conditional	✓
Pred6	Test with addition	✗
Pred7	Test with multiple variables	✗

formulated the certification conditions of Denning’s analysis [18] as a simple type system for a deterministic language: basically, a formal system of type inference rules for making judgments about programs. More generally, type-based approaches are designed such that well-typed programs do not leak secrets. A type is inductively associated at compile-time with program statements in such a way that any statement showing a potential low disclosing secrets is rejected. Type systems that enforce secure information flow have been designed for various languages and they have been used in different applications. Some of these approaches are, for example, applied to specific programs, e.g., written in VHDL [44], where the analysis of information flow is closely related to the context. Moreover, the secure information flow problem was also handled in different situations, for example with multi-threaded programs [42] or with programs that employ explicit cryptographic operations [21, 3].

A different approach is the use of standard control flow analysis to detect information leakage, e.g., [9, 29, 30]. The idea, of this technique, is to conservatively find the program paths through which data may flow. Generally, the data flow analysis approach to secure information flow as a translation from a given program that captures and facilitates reasoning about the possible flows. For example, Leino and Joshi [29] showed an application based on semantics, deriving a first-order predicate whose validity implies that an attacker cannot

Table 20: Jif case studies

Name	Description	fa
A	Simple explicit flow test	✓
Account	Simple explicit flow test	✓
ConditionalLeak	Explicit flow in if statement	✓
Do	Implicit flow in the loop	✓
Do2	Implicit flow if and loop	✓
Do3	Implicit flow loop and if	✓
Do4	Implicit flow loop and if	✓
Do5	Implicit flow loop and if	✓
If1	Simple implicit flow	✓
Implicit	Simple implicit flow	✓

deduce any secure information from observing the public inputs, outputs and termination behavior of the program.

The use of abstract interpretation in language-based security is not new, even though there aren't many works that use the lattice of abstract interpretations for evaluating the security of programs (for example [49]).

Probably, the main work about information flow analysis by abstract interpretation was done by Giacobazzi and Mastroeni [22] that generalizes the notion of non-interference making it parametric relatively to what an attacker can observe, and using it to model attackers as abstractions. A program semantics was characterized as an abstract interpretation of its maximal trace semantics in the corresponding transition system. The authors gave a method for checking abstract non-interference and they proved that checking abstract non-interference is a standard static program analysis problem. This method allows both to compare attackers and program secrecy by comparing the corresponding abstractions in the lattice of abstract interpretations, and to design automatic program certification tools for language-based security.

There are not so many implementations of secure information flow. In early 2000, some works began the control of sensitive information in realistic languages [7, 37]. Jif [4] and Flow CAML [40] are, as far as we know, the two main implementations about information flow analysis. Notice that, in the last years other language-based tools are developed for some specific language, e.g., Fabric [32] for distributed computing, the LIO library in haskell [43] and FlowFox([16]) a tool for JavaScript.

According to [41], it seems be helpful to distinguish between two different application scenarios: *developing secure software* and *stopping malicious software*. The first scenario is based on to secure information flow analysis to help the development of software that satisfies some security properties. In this case, the analysis serves as a program development tool. The static analysis tool would alert the programmer to potential leaks and the developer could rewriting the code as necessary. An example of this scenario can be found in [4], where Askarov and Sabelfeld discusses the implementation of a "mental poker" protocol in Jif. The second scenario, instead, the secure information flow analysis is used as a

kind of filter to stop malicious software. In this case, we might imagine analyzing a piece of untrusted code before executing it, with the goal of guaranteeing its safety. This is much more challenging than first scenario: probably we would not have access to the source code and we would need to analyze the binary code. Analyzing binaries is more difficult than analyzing source code and has not received much attention in the literature (a Java bytecodes analysis is performed, for instance, by Barthe and Rezk in [8]).

Given this overall context, the approach adopted in *Sails* is quite different from existing tools that deal with information flow analysis. *Jif*, for example, is a security-typed programming language that extends Java with support for information flow and access control, enforced at compile time and it is an ad hoc analysis that requires to annotate the code with some type information. If on the one hand *Jif* is more efficient than *Sails*, on the other hand *Sails* does not require any manual annotation, and it takes all advantages of compositional analyzers (e.g., we can combine *Sails* with a TVLA-based heap abstraction).

Our approach does not require to change the programming language, since it infers the flow of information directly on the original program, and it asks what are the private data that have not to be leaked to the user during the analysis execution.

## 8 Conclusions

In this paper we presented an information flow analysis through abstract interpretation based on a new domain that combines a variable dependency analysis and a numerical domain. We then introduced *Sails* that applies and implements this analysis on object-oriented programs. *Sails* is an extension of *Sample*, therefore it is modular with respect to the heap abstraction, and it can verify noninterference over recursive data structures using simple and efficient heap analyses. The experimental results underline the effectiveness of the analysis, since *Sails* is in position to analyze several benchmarks in few milliseconds per program without producing false alarms in more than 90% of the programs. Moreover, our tool does not require to modify the original language, since it works with mainstream languages like Java, and it does not require any manual annotation.

**Acknowledgments** This work has been partially supported by CINI Cybersecurity National Laboratory within the project "FiliereSicura: Securing the Supply Chain of Domestic Critical Infrastructures from Cyber Attacks" funded by CISCO Systems Inc. and Leonardo SpA, and by MIUR-MAE within the Project "Formal Specification for Secured Software System", under the Indo-Italian Executive Programme of Cooperation in Scientific and Technological Cooperation Project number IN17MO07.

## References

1. Henrik Reif Andersen. An introduction to binary decision diagrams. Technical report, Course Notes on the WWW, 1997.

2. Tania Armstrong, Kim Marriott, Peter Schachte, and Harald Søndergaard. Two classes of boolean functions for dependency analysis. *Sci. Comput. Program.*, 31:3–45, May 1998.
3. Aslan Askarov, Daniel Hedin, and Andrei Sabelfeld. Cryptographically-masked flows. *Theor. Comput. Sci.*, 402:82–101, July 2008.
4. Aslan Askarov and Andrei Sabelfeld. Security-typed languages for implementation of cryptographic protocols: A case study. In *ESORICS*, Lecture Notes in Computer Science, pages 197–221, 2005.
5. Roberto Bagnara, Patricia M. Hill, and Enea Zaffanella. The parma polyhedra library: Toward a complete set of numerical abstractions for the analysis and verification of hardware and software systems. *Sci. Comput. Program.*, 72:3–21, June 2008.
6. Roberto Bagnara, Patricia M. Hill, and Enea Zaffanella. Applications of polyhedral computations to the analysis and verification of hardware and software systems. *Theor. Comput. Sci.*, 410:4672–4691, November 2009.
7. Anindya Banerjee and David A. Naumann. Secure information flow and pointer confinement in a java-like language. In *Proceedings of the 15th IEEE workshop on Computer Security Foundations, CSFW '02*, Washington, DC, USA, 2002. IEEE Computer Society.
8. Gilles Barthe and Tamara Rezk. Non-interference for a jvm-like language. In *Proceedings of the 2005 ACM SIGPLAN international workshop on Types in languages design and implementation, TLDI '05*, pages 103–112, New York, NY, USA, 2005. ACM.
9. Chiara Bodei, Pierpaolo Degano, Flemming Nielson, and Hanne Riis Nielson. Static analysis for secrecy and non-interference in networks of processes. In *Proceedings of the 6th International Conference on Parallel Computing Technologies, PaCT '01*, pages 27–41, London, UK, 2001. Springer-Verlag.
10. James Clause, Wanchun Li, and Alessandro Orso. Dytan: a generic dynamic taint analysis framework. In *Proceedings of the 2007 international symposium on Software testing and analysis, ISSTA '07*, pages 196–206, New York, NY, USA, 2007. ACM.
11. Agostino Cortesi, Gilberto Filé, and William H. Winsborough. Prop revisited: Propositional formula as abstract domain for groundness analysis. In *LICS*, pages 322–327, 1991.
12. Giulia Costantini, Pietro Ferrara, and Agostino Cortesi. Static analysis of string values. In *ICFEM*, Lecture Notes in Computer Science. Springer, 2011.
13. Patrick Cousot. The Calculational Design of a Generic Abstract Interpreter. In M. Broy and R. Steinbrüggen, editors, *Calculational System Design*. NATO ASI Series F. IOS Press, Amsterdam, 1999.
14. Patrick Cousot and Radhia Cousot. Systematic design of program analysis frameworks. In *Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages, POPL '79*, pages 269–282, New York, NY, USA, 1979. ACM.
15. Patrick Cousot and Nicolas Halbwachs. Automatic discovery of linear restraints among variables of a program. In *Proceedings of the 5th ACM SIGACT-SIGPLAN symposium on Principles of programming languages, POPL '78*, pages 84–96, New York, NY, USA, 1978. ACM.
16. Willem De Groef, Dominique Devriese, Nick Nikiforakis, and Frank Piessens. FlowFox: a web browser with flexible and precise information flow control. In *Proceedings of the 19th ACM Conference on Computer and Communications Security (CCS 2012)*. ACM, 2012.
17. Dorothy E. Denning. A lattice model of secure information flow. *Commun. ACM*, 19:236–243, May 1976.
18. Dorothy E. Denning and Peter J. Denning. Certification of programs for secure information flow. *Commun. ACM*, 20:504–513, July 1977.

19. P. Ferrara. Static type analysis of pattern matching by abstract interpretation. In *Formal Techniques for Distributed Systems (FMOODS/FORTE)*, volume 6117 of *Lecture Notes in Computer Science*, pages 186–200. Springer-Verlag, 2010.
20. Pietro Ferrara. A fast and precise alias analysis for data race detection. In *Proceedings of the Third Workshop on Bytecode Semantics, Verification, Analysis and Transformation (Bytecode'08)*, *Electronic Notes in Theoretical Computer Science*. Elsevier, April 2008.
21. Riccardo Focardi and Matteo Centenaro. Information flow security of multi-threaded distributed programs. In *Proceedings of the third ACM SIGPLAN workshop on Programming languages and analysis for security, PLAS '08*, pages 113–124, New York, NY, USA, 2008. ACM.
22. Roberto Giacobazzi and Isabella Mastroeni. Abstract non-interference: parameterizing non-interference by abstract interpretation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on Principles of programming languages, POPL '04*, pages 186–197, New York, NY, USA, 2004. ACM.
23. Raju Halder and Agostino Cortesi. Abstract interpretation of database query languages. *Computer Languages, Systems & Structures*, 38:123–157, 2012.
24. Raju Halder and Agostino Cortesi. Abstract program slicing of database query languages. In *Proc. of the 28th Annual ACM Symposium on Applied Computing*, pages 838–845, Coimbra, Portugal, 2013. ACM Press.
25. Raju Halder, Matteo Zanioli, and Agostino Cortesi. Information leakage analysis of database query languages. In *Proc. of the 29th Annual ACM Symposium on Applied Computing*, pages 813–820, Gyeongju, Korea, March 24–28 2014. ACM Press.
26. Matthew Hennessy. *The semantics of programming languages: an elementary introduction using structural operational semantics*. John Wiley & Sons, Inc., New York, NY, USA, 1990.
27. B. Jeannet. Convex polyhedra library, March 2002. Documentation of the “New Polka” library available at <http://www.irisa.fr/prive/Bertrand.Jeannet/newpolka.html>.
28. B. Jeannet and A. Miné. Apron: A library of numerical abstract domains for static analysis. In *Proc. of the 21th Int. Conf. on Computer Aided Verification (CAV 2009)*, volume 5643 of *Lecture Notes in Computer Science*, pages 661–667. Springer, June 2009.
29. Rajeev Joshi and K. Rustan M. Leino. A semantic approach to secure information flow. *Sci. Comput. Program.*, 37:113–138, May 2000.
30. Peeter Laud. Semantics and program analysis of computationally secure information flow. In *Proceedings of the 10th European Symposium on Programming Languages and Systems, ESOP '01*, pages 77–91, London, UK, 2001. Springer-Verlag.
31. Peng Li and Steve Zdancewic. Arrows for secure information flow. *Theor. Comput. Sci.*, 411:1974–1994, April 2010.
32. Jed Liu, Michael D. George, K. Vikram, Xin Qi, Lucas Waye, and Andrew C. Myers. Fabric: a platform for secure distributed computation and storage. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles, SOSP '09*, pages 321–334, New York, NY, USA, 2009. ACM.
33. Yin Liu and Ana Milanova. Static information flow analysis with handling of implicit flows and a study on effects of implicit flows vs explicit flows. In *Proceedings of the 2010 14th European Conference on Software Maintenance and Reengineering, CSMR '10*, pages 146–155, Washington, DC, USA, 2010. IEEE Computer Society.
34. Antoine Miné. A new numerical abstract domain based on difference-bound matrices. In *Proceedings of the Second Symposium on Programs as Data Objects, PADO '01*, pages 155–172, London, UK, 2001. Springer-Verlag.
35. Antoine Miné. The octagon abstract domain. In *Proc. of the Workshop on Analysis, Slicing, and Transformation (AST'01)*, IEEE, pages 310–319. IEEE CS Press, October 2001.

36. A. C. Myers, L. Zheng, S. Zdancewic, S. Chong, and N. Nystrom. Jif: Java information flow. software release., July 2001-2004.
37. François Pottier and Vincent Simonet. Information flow inference for ml. *ACM Trans. Program. Lang. Syst.*, 25:117–158, January 2003.
38. A. Sabelfeld and A. C. Myers. Language-based information-flow security. *Selected Areas in Communications, IEEE Journal on*, 21(1):5–19, 2003.
39. Mooly Sagiv, Thomas Reps, and Reinhard Wilhelm. Parametric shape analysis via 3-valued logic. *ACM Trans. Program. Lang. Syst.*, 24:217–298, May 2002.
40. Vincent Simonet. The Flow Caml System: documentation and user’s manual. Technical Report 0282, Institut National de Recherche en Informatique et en Automatique (INRIA), July 2003. ©INRIA.
41. Geoffrey Smith. Principles of Secure Information Flow Analysis. In *Malware Detection*, pages 297–307, 2007.
42. Geoffrey Smith and Dennis Volpano. Secure information flow in a multi-threaded imperative language. In *Proceedings of the 25th ACM SIGPLAN-SIGACT symposium on Principles of programming languages, POPL ’98*, pages 355–364, New York, NY, USA, 1998. ACM.
43. Deian Stefan, Alejandro Russo, John C. Mitchell, and David Mazières. Flexible dynamic information flow control in haskell. *SIGPLAN Not.*, 46(12):95–106, September 2011.
44. T. K. Tolstrup, F. Nielson, and H. Riis Nielson. Information flow analysis for VHDL. In *Parallel Computing Technologies*, Lecture Notes in Computer Science. Springer, 2005.
45. Stanford University. Stanford SecuriBench Micro. <http://suif.stanford.edu/~livshits/work/securibench-micro/>.
46. Dennis Volpano, Cynthia Irvine, and Geoffrey Smith. A sound type system for secure flow analysis. *J. Comput. Secur.*, 4:167–187, January 1996.
47. Matteo Zanioli and Agostino Cortesi. Information leakage analysis by abstract interpretation. In *Proceedings of the 37th international conference on Current trends in theory and practice of computer science (SOFSEM)*, volume 6543 of *Lecture Notes in Computer Science*, pages 545–557. Springer-Verlag, 2011.
48. Matteo Zanioli, Pietro Ferrara, and Agostino Cortesi. Sails: Static analysis of information leakage with sample. In *Proceedings of the 2012 ACM symposium on Applied Computing*, pages 1308–1313. ACM Press, 2012.
49. Mirko Zanotti. Security typings by abstract interpretation. In *Proceedings of the 9th International Symposium on Static Analysis, SAS ’02*, pages 360–375, London, UK, 2002. Springer-Verlag.