

Formal methods : implementing a secure-type system for an imperative programming language

Youssef MILED

June 2021

Abstract

Connected devices are becoming ever more prone to cyberattacks, emphasizing the pressing need for secure software. Formal methods represent an approach based on mathematical principles to guarantee program correctness, and language-based security (LBS) offers a certain technique for information flow control. This paper studies security of information flow in a simple imperative language with constructs for expressions, assignments, conditionals, loops and function calls. We expand on Geoffrey Smith's type system to cover functions, and we classify information based on confidentiality (public vs. private). Typing rules establish constraints that a code must meet to be assigned a specific type, ensuring that well-typed commands adhere to specific security constraints, while subtyping rules define relationships between types. Our research revolves around the development and integration of a type system that will enhance software security in a simple programming language. With OCaml as the basis of this study, we implement the type system to enforce the security measures. The central interest lies in assessing how effective the type system is at preventing illegal data flow and in enhancing confidentiality and integrity restrictions, by studying and proving the non-interference property to the extended type system. This is done by also studying equivalence between memories. In summary, this approach joins both theoretical and practical work to develop stronger software security solutions.

Contents

1	Introduction	2
1.1	Motivation	2
1.2	Background	2
1.3	Objective	2
1.4	Contributions	2
2	Modelisation	3
2.1	Language	3
2.2	Semantics	4
3	Type system	4
3.1	Overview	4
3.2	Properties	7
4	Implementation	9
5	Conclusion	10
5.1	Summary	10
5.2	Future directions	10
6	Appendix	12

1 Introduction

1.1 Motivation

The vulnerability of connected objects to cyberattacks in the connected city is an issue that prompts search for effective methods to tackle problems of illegally obtaining access to the system or theft of confidential data, by securing the computer programs involved in these processes.

These systems can be developed using formal methods which will ensure they meet certain specifications thus improving their security.

1.2 Background

Formal methods play a significant role in securing computer systems. Beyond testing, they enable the assimilation of a program into a mathematical definition using certain semantics, thus employing logical deduction to establish true properties for all possible program inputs without having to test.

One can choose to reason about abstract models such as program graphs (PG). For instance, program verification involves associating predicates with each node of the PG. This approach remains challenging to automate insofar as predicates are difficult to express in machine. Another method involves the language-based security (LBS) approach ensures the security of a program in a given programming language.

Indeed, this approach is based notably on the study of the information flow allowed during program execution. Various security levels are then distinguished for each input or output involved in the program: this generally translates into sets {reliable, dubious} for integrity, {private, public} or also {unclassified, classified, secret, top secret} for confidentiality.

Controlling the information flow means determining whether a program can allow a sequence of information from a variable x , with low confidentiality level, to a variable y , with high confidentiality level, either through a direct (explicit) flow or through a conditional branch and other mechanisms (implicit flow). Information Flow Control (IFC) has already been applied to certain programming languages such as Caml (Flow Caml).

Besides, Nevin Heintze and Jon G. Riecke designed the SLam calculus (Secure Lambda Calculus), which is a typed lambda calculus containing, in addition to the usual types, security information such as reader, indirect reader (aiming to specify confidentiality), as well as creator and indirect creator (aiming to specify integrity).

Therefore using a programming language logic (semantics, rules) for security purposes allows statically to verify programs before their execution.

1.3 Objective

This work revolves around exploring a specific type system customised for a simple programming language, in order to determine its contribution to the security of computer programs. Furthermore, the ultimate goal consists in translating theoretical insights into practical implementation by implementing the identified approach within the OCaml programming environment.

A proper combination of the theoretical grounds and the practical deployment allows generating the results that could be further exploited in secure coding initiatives and security for software systems.

We use formal methods to guarantee a certain level of safety; for instance the absence of unintended data transfer within the program.

1.4 Contributions

The main contributions of this paper are:

- We develop a certain type system for secure information flow by extending the use of functions and typing rules to Geoffrey Smith's grammar, type system and semantics.
- We establish the validity of the type system by proving the non-interference theorem for the added functionalities.
- We make the case for implementing a representation of trees for secure type proofs.

2 Modelisation

2.1 Language

The language used in this work is a simple imperative language that contains the key aspects a programming language must have, i.e assignments, conditional branches, while loops and binary operations between expressions.

The language that we work on is defined by the following grammar:

$$\begin{aligned}
e &::= n \mid x \mid e_1 \odot e_2 \mid f(e) \\
c &::= x := e \mid \mathbf{skip} \mid \mathbf{if} \ e \ \mathbf{then} \ c_1 \ \mathbf{else} \ c_2 \mid \\
&\quad \mathbf{while} \ e \ \mathbf{do} \ c \mid c_1; \ c_2 \mid f(e) \mid \mathbf{return} \ e \mid \\
f &::= fun \ (\lambda.x) \ \{c\} \\
f &::= fun \ (\lambda.x) \{c\} \\
p &::= e \mid c
\end{aligned}$$

The grammar we present includes expressions (e), commands (c), and phrases (p). It defines a language that contains basic expressions, while loops, conditional statements (if-else), and the ability to sequence commands. This framework provides the essential tools needed to build simple programs, allowing for both basic calculations and control flow structures.

Each category (expression, command, phrase) is essential to the language's structure, serving distinct roles within the language's syntax and semantics. We also add functions, which are useful to reuse codes.

EXPRESSIONS:

Expressions (e) are the basic computational entities. They are defined recursively as follows:

Numerals (n): represent for integer constants.

Variables (x): represent symbols storing values.

Binary Operations ($e_1 \odot e_2$): combinatory structures which enable the construction of expressions. $\odot \in \{+, -, \times, \div, =, \neq, <, \geq, \wedge, \vee\}$.

COMMANDS:

Commands (c) contain the execution or control structures within the language:

Assignment ($x := e$): The variable x is assigned the value obtained by evaluating expression e .

No-operation (**skip**): This command does not perform any action; it is useful in complex sequences for expressing semantic no-ops. We will see later how this command is crucial so that the type system offers secure information flow.

Conditional execution (**if** e **then** c_1 **else** c_2): Directs the flow of execution through the evaluation of expression e . If e evaluates to a true value, command c_1 is executed; otherwise, c_2 is executed. The else branch is important for the type system in order to offer secure information flow.

Iteration (**while** e **do** c): Implements the looping behavior in which command c will be repeated as long as the expression e is evaluated to true.

Sequence ($c_1; c_2$) is for the sequential execution of two commands.

Functions ($fun(\lambda.x)\{c\}$) enable to reuse code and to organize it into small pieces or modules. Functions take in inputs which are specific types and return an output of another type. We reserve the notation $\lambda.x$ for an internal variable within a function.

PHRASES:

Phrases (p) are the top-level syntactic units of the language, which can be either an expression or a command. In this system, code is therefore defined by a certain phrase p .

An example of a code in this grammar could be a phrase defined by :

$$a := 5; \ b := 2; \ \mathbf{while} \ a \ \mathbf{do} \ a := a - b$$

The design of this grammar supports the development of a structured programming language capable of expressing a wide range of computational logic through a compact and coherent syntactic framework.

This approach does not help in clearly defining computational semantics, a key factor to understand the theoretical foundations and practical applications of the language.

2.2 Semantics

In this language, programs are executed based on a specific memory denoted as μ , which consists of a set of pairs $\{(e_i, n_i)\}$, where e_i represents expressions and n_i represents integers.

The semantics of commands are defined through a sequential transition relation denoted as \longrightarrow on configurations.

Definition 2.1. A configuration C is defined by a couple (c, μ) or μ where c is a command and μ is a memory.

We extend Geoffrey Smith's semantics to include functions here ((FUNC) semantic).

(UPDATE)	$\frac{x \in \text{dom}(\mu)}{(x := e, \mu) \longrightarrow \mu[x := \mu(e)]}$
(NO-OP)	$(\text{skip}, \mu) \longrightarrow \mu$
(BRANCH)	$\frac{\mu(e) \neq 0}{(\text{if } e \text{ then } c_1 \text{ else } c_2, \mu) \longrightarrow (c_1, \mu)}$
	$\frac{\mu(e) = 0}{(\text{if } e \text{ then } c_1 \text{ else } c_2, \mu) \longrightarrow (c_2, \mu)}$
(LOOP)	$\frac{\mu(e) = 0}{(\text{while } e \text{ do } c, \mu) \longrightarrow \mu}$
	$\frac{\mu(e) \neq 0}{(\text{while } e \text{ do } c, \mu) \longrightarrow (c; \text{while } e \text{ do } c, \mu)}$
(SEQUENCE)	$\frac{(c_1, \mu) \longrightarrow \mu'}{(c_1; c_2, \mu) \longrightarrow (c_2, \mu')}$
	$\frac{(c_1, \mu) \longrightarrow (c'_1, \mu')}{(c_1; c_2, \mu) \longrightarrow (c'_1; c_2, \mu')}$
(RET)	$(\text{return } e, \mu) \longrightarrow \mu$
(FUNC)	$(f ::= \text{fun } (\lambda.x) \{c\}, \mu) \rightarrow \mu$
	$\frac{\begin{array}{l} \text{dom}(\mu'') = \text{dom}(\mu') \setminus \{\lambda.x\} \\ \mu''(e) = \mu'(\lambda.x) \\ (c, \mu \cup \{(\lambda.x, \mu(e))\}) \longrightarrow^* \mu' \end{array}}{(f(e)\{c\}, \mu) \longrightarrow \mu''}$

Table 1: Extended semantics

Here, for the sake of simplicity, we assume that function calls are evaluated atomically. To ensure this we can use a mutex.

We explain the semantic for functions as follows: in order for a function call $f(e)\{c\}$ to modify a memory μ into μ'' , executing the command c with the memory μ extended with the internal variable $\lambda.x$ should result in a memory μ' and we also require $\mu''(e) = \mu'(\lambda.x)$ since $e \in \text{dom}(\mu'')$ which is the memory outside the function. Naturally, the final memory μ'' does not include the internal variable $\lambda.x$.

3 Type system

3.1 Overview

We revisit Georges Smith type system for our language and add other features such as functions. This revisited system will be proven to verify information flow security (non-interference property, see 1). The

types used in this system are:

$$\begin{aligned}\tau &::= L \mid H \\ \rho &::= \tau \mid \tau \text{ cmd } n \mid \tau_1 \text{ cmd } \tau_2 \mid \tau_1 \rightarrow \tau_2\end{aligned}$$

Indeed, there are two main classes in Geoffrey Smith's (GS) type system: L (for low) which characterizes a public information, that does not need to be protected, and H (for high) which is made for private information that should be secured. Besides, commands c which are typed $\tau_1 \text{ cmd } \tau_2$ assign only to variables of type τ_1 and whose execution time depends on variables of type τ_2 . The same thing goes for commands of type $\tau \text{ cmd } n$, which execute in n steps.

We introduce a new type constructor for functions, that GS's type system does not include. A function of type $\tau_1 \rightarrow \tau_2$ is a function that has an argument of type τ_1 and returns a value of type τ_2 .

Therefore, programs are constructed following the revisited rules from Geoffrey Smith typing rules (TABLE 2). We added to these some rules involving functions.

(R-VAL)	$\frac{\gamma(x) = \tau \text{ var}}{\gamma \vdash x : \tau}$
(INT)	$\gamma \vdash n : L$
(BINOP)	$\frac{\gamma \vdash e_1 : \tau_1, \gamma \vdash e_2 : \tau_2}{\gamma \vdash e_1 \odot e_2 : \tau}$
(ASSIGN)	$\frac{\gamma(x) = e : \tau \text{ var}, \gamma \vdash e : \tau}{\gamma \vdash x := e : \tau \text{ cmd } 1}$
(SKIP)	$\gamma \vdash \mathbf{skip} : H \text{ cmd } 1$
(IF)	$\frac{\begin{array}{c} \gamma \vdash e : \tau \\ \gamma \vdash c_1 : \tau \text{ cmd } n \\ \gamma \vdash c_2 : \tau \text{ cmd } n \end{array}}{\gamma \vdash \mathbf{if } e \mathbf{ then } c_1 \mathbf{ else } c_2 : \tau \text{ cmd } n+1}$
	$\frac{\begin{array}{c} \gamma \vdash e : L \\ \gamma \vdash c_1 : \tau_1 \text{ cmd } \tau_1 \\ \gamma \vdash c_2 : \tau_1 \text{ cmd } \tau_2 \end{array}}{\gamma \vdash \mathbf{if } e \mathbf{ then } c_1 \mathbf{ else } c_2 : \tau_1 \text{ cmd } \tau_2}$
	$\frac{\begin{array}{c} \gamma \vdash e : H \\ \gamma \vdash c_1 : H \text{ cmd } H \\ \gamma \vdash c_2 : H \text{ cmd } H \end{array}}{\gamma \vdash \mathbf{if } e \mathbf{ then } c_1 \mathbf{ else } c_2 : H \text{ cmd } H}$
(WHILE)	$\frac{\begin{array}{c} \gamma \vdash e : L \\ \tau_2 \subseteq \tau_1 \\ \gamma \vdash c : \tau_1 \text{ cmd } \tau_2 \end{array}}{\gamma \vdash \mathbf{while } e \mathbf{ do } c : \tau_1 \text{ cmd } \tau_2}$
	$\frac{\begin{array}{c} \gamma \vdash e : H \\ \gamma \vdash c : H \text{ cmd } H \end{array}}{\gamma \vdash \mathbf{while } e \mathbf{ do } c : \tau_1 \text{ cmd } \tau_2}$
(COMPOSE)	$\frac{\begin{array}{c} \gamma \vdash c_1 : \tau_1 \text{ cmd } L \\ \gamma \vdash c_2 : \tau_1 \text{ cmd } \tau_2 \end{array}}{\gamma \vdash c_1; c_2 : \tau_1 \text{ cmd } \tau_2}$
	$\frac{\begin{array}{c} \gamma \vdash c_1 : \tau \text{ cmd } H \\ \gamma \vdash c_2 : H \text{ cmd } H \end{array}}{\gamma \vdash c_1; c_2 : \tau \text{ cmd } H}$

Table 2: Geoffrey Smith typing rules

We add a rule for the **return** command.

$$(\text{RET}) \quad \frac{\gamma \vdash e : \tau_1}{\gamma \vdash \text{return } e : \tau_1 \text{ cmd } 1}$$

Considering a function f defined by $f ::= \text{fun}(\lambda.x)\{c; \text{return } e\}$:

$$(\text{FUNC}) \quad \frac{\gamma, \lambda.x : \tau_1 \vdash y : \tau_2 \quad \gamma, \lambda.x : \tau_1 \vdash c : \tau_1 \text{ cmd } n}{\gamma \vdash f(\lambda.x)\{c; \text{return } y\} : \tau_1 \rightarrow \tau_2}$$

If f is defined by $f ::= \text{fun}(\lambda.x)\{c\}$:

$$(\text{VOID-FUNC}) \quad \frac{\gamma, \lambda.x : \tau_1 \vdash c : \tau_1 \text{ cmd } n}{\gamma \vdash f(\lambda.x)\{c\} : \tau_1 \rightarrow H}$$

In both cases, the rule to call a function is as follows :

$$(\text{CALL}) \quad \frac{\gamma \vdash e : \tau_1 \quad \gamma \vdash f(\lambda.x)\{c\} : \tau_1 \rightarrow \tau_2}{\gamma \vdash f(e)\{c\} : \tau_2}$$

We presume that upon calling a function f with a given expression e , the command $\lambda.x := e$ is automatically executed.

For a definition of a function, nothing should impact the code and the information flow, therefore we consider the following rule:

$$(\text{FUNC-DEF}) \quad \gamma \vdash f ::= \text{fun } (\lambda.x)\{c\} : L \text{ cmd } H$$

Indeed, for any command type $\tau_1 \text{ cmd } \tau_2$, we have $L \text{ cmd } H \subseteq \tau_1 \text{ cmd } \tau_2$ (see subtyping rules 3).

Note that γ , the typing environment, contains the types of the variables that are manipulated for a certain code.

These new typing rules involve two categories of functions :

- functions that return an expression. For instance,

$$f(\lambda.x)\{\lambda.y := x + 1; \text{return } \lambda.y\}$$

is in this category.

- functions that do not return anything and simply apply certain commands. For instance,

$$f(\lambda.x)\{\lambda.x := \lambda.x + 1\}$$

is a void function.

In order to use these rules in formal proofs, we need to define subtyping rules which define an order for the types that we defined. For this matter, we extend Geoffrey Smith's subtyping rules by introducing rules that establish an order among the new added types.

(BASE)	$L \subseteq H$
(CMD)	$\frac{\tau_1 \subseteq \tau'_1, \tau'_2 \subseteq \tau_2}{\tau_1 \text{ cmd } \tau_2 \subseteq \tau'_1 \text{ cmd } \tau'_2}$ $\frac{\tau' \subseteq \tau}{\tau \text{ cmd } n \subseteq \tau' \text{ cmd } n}$ $\tau \text{ cmd } n \subseteq \tau \text{ cmd } L$
(REFLEX)	$\rho \subseteq \rho$
(TRANS)	$\frac{\rho_1 \subseteq \rho_2, \rho_2 \subseteq \rho_3}{\rho_1 \subseteq \rho_3}$
(SUBSUMP)	$\frac{\gamma \vdash p : \rho_1, \rho_1 \subseteq \rho_2}{\gamma \vdash p : \rho_2}$

Table 3: Geoffrey Smith subtyping rules

Here are the additional subtyping rules:

(FUNC)	$\tau_1 \rightarrow \tau_2 \subseteq \tau_1 \text{ cmd } 1$
	$\frac{\tau'_1 \subseteq \tau_1, \tau_2 \subseteq \tau'_2}{\tau_1 \rightarrow \tau_2 \subseteq \tau'_1 \rightarrow \tau'_2}$

Table 4: New subtyping rules

The last rule ((FUN)) can be justified as follows: A function that takes as an argument $x : \tau_1$ can also be viewed as a general case of a function that takes an argument of type $\tau'_1 \subseteq \tau_1$ and returns an expression of type τ'_2 since it is designed to return an expression of type $\tau_2 \subseteq \tau'_2$.

3.2 Properties

Before delving into the property, let's recall some definitions.

Definition 3.1. A command c is well typed with respect to a typing environment γ iff $\exists \rho. \gamma \vdash c : \rho$.

Definition 3.2. 2 memories μ and ν are equivalent with respect to γ iff μ and ν have the same domain and agree on all L variables.

Definition 3.3. 2 commands c and d are equivalent with respect to γ iff : c and d are well typed with respect to γ and one of the following conditions is verified :

- $c = d$
- $c : H \text{ cmd } \tau$ and $d : H \text{ cmd } \tau$
- $c = c_1; c_2; \dots; c_k$ and $d = d_1; c_2; \dots; c_k$ with c_1 and $d_1 : H \text{ cmd } n$

This relation is denoted by \sim_γ .

Definition 3.4. 2 Configurations C and D are equivalent with respect to γ iff

- $C = (c, \mu), D = (d, \nu)$ with $c \sim_\gamma d$ and $\mu \sim_\gamma \nu$
- $C = (c, \mu), D = \nu$ with $c : H \text{ cmd } \tau$ and $\mu \sim_\gamma \nu$.
- $C = \mu$ and $D = \nu$ with $\mu \sim_\gamma \nu$.

Lemma 1. Given two memories μ and ν and a typing environment γ , if $\mu \sim_\gamma \nu$ then

$$\mu \cup \{(\lambda.x, \mu(e))\} \sim_{\gamma'} \nu \cup \{(\lambda.x, \nu(e))\}$$

where $\gamma' = \gamma \cup \{(\lambda.x, \tau \text{ var})\}$.

Proof. Let's consider two memories μ and ν and γ a typing environment. We suppose $\mu \sim_\gamma \nu$. Let $\gamma' = \gamma \cup \{(\lambda.x, \tau \text{ var})\}$. To prove that $\mu \sim_{\gamma'} \nu$ we should prove that they agree on all L variables in γ' . Let $(y, L \text{ var}) \in \gamma'$. If $(y, L \text{ var}) \in \gamma$ then $\mu(y) = \nu(y)$ because $\mu \sim_\gamma \nu$. Otherwise, $(y, L \text{ var}) = (\lambda.x, \tau \text{ var})$ so $\tau = L$. Thus, $\gamma \vdash e : L$ according to the rule ASSIGN in 2. So $\mu(e) = \nu(e)$ and $\mu(y) = \nu(y)$ because we have $\lambda.x := e$ when calling $f(e)$. In both cases $\mu(y) = \nu(y)$. Therefore, $\mu \sim_{\gamma'} \nu$. \square

Lemma 2. *Given two memories μ and ν and a typing environment γ , if $\mu \sim_\gamma \nu$ and $\gamma' \subseteq \gamma$ then $\mu \sim_{\gamma'} \nu$.*

Proof. If μ and ν agree on all L variables in γ , they also agree on all L variables in γ' since $\gamma' \subseteq \gamma$. \square

The type system that we developed here, based on Geoffrey Smith's and some added functionalities verifies the property of non-interference.

Theorem 1. *(non-interference).*

$$\begin{aligned} & \forall \gamma \in F(I, \{L \text{ var}, H \text{ var}\}), \forall c \in \text{commands}, \forall \mu, \nu \in F(I, \mathbb{N}), \\ & (\exists \rho. \gamma \vdash c : \rho \wedge \mu \sim_\gamma \nu \wedge (c, \mu) \longrightarrow C' \wedge (c, \nu) \longrightarrow D') \implies (C' \sim_\gamma D') \end{aligned}$$

Proof. Let c be a well typed command with respect to γ , i.e $\exists \rho. \gamma \vdash c : \rho$. Suppose also $\mu \sim_\gamma \nu$ and $(c, \mu) \longrightarrow C'$ and $(c, \nu) \longrightarrow D'$. For the general case, we consider $c = c_1; c_2; \dots; c_k$ and we consider in turn each of the possible forms of c_1 . The theorem is proven in Geoffrey Smith's paper if c_1 is of one of these forms :

- $x := e$
- **skip**
- **if** e' **then** $c_{1,1}$ **else** $c_{1,2}$
- **while** e **do** $c_{1,1}$

If c_1 is of the form $f := \text{fun } (\lambda.x) \{d\}$, i.e c_1 defines a function, then according to the first semantic rule (FUNC) 1, we can deduce that $C' = (c_2; \dots; c_k, \mu)$ and $D' = (c_2; \dots; c_k, \nu)$ and $C' \sim_\gamma D'$ holds true because of the equivalence between μ and ν .

We now consider c_1 as $f(e)$ for some expression e and some void function $f(\lambda.x)\{d\}$. We suppose that d is a command without any void function call, therefore $d = d_1; d_2; \dots; d_k$ with d_i defined as **if** e' **then** $d_{i,1}$ **else** $d_{i,2}$ **while** e **do** $d_{i,1}$, or **skip**. We can always reduce to this case since when demonstrating the validity of this specific case, we can then build our argument for all other possible cases iteratively.

Since $(c, \mu) \longrightarrow C'$, by the semantic (SEQUENCE) we have $C' = (c_2; \dots; c_k, \mu')$ and similarly, D' is $(c_2; \dots; c_k, \nu')$, where μ' (respectively ν') is the modified memory μ (respectively ν) after applying the void function f according to the following semantic:

$$\frac{\begin{array}{l} \text{dom}(\mu') = \text{dom}(\mu'') \setminus \{\lambda.x\} \\ \mu'(e) = \mu''(\lambda.x) \\ (d, \mu \cup \{(\lambda.x, \mu(e))\}) \longrightarrow^* \mu'' \end{array}}{(f(e)\{d\}, \mu) \longrightarrow \mu'}$$

Thus, it follows that $(d, \mu \cup \{(\lambda.x, \mu(e))\}) \longrightarrow^* \mu''$, which means that executing d with the memory $\mu \cup \{(\lambda.x, \mu(e))\}$ leads to the memory μ'' . Similarly, $(d, \nu \cup \{(\lambda.x, \nu(e))\}) \longrightarrow^* \nu''$. Besides, note that the typing environment corresponding to $\mu \cup \{(\lambda.x, \mu(e))\}$ and $\nu \cup \{(\lambda.x, \nu(e))\}$ is $\alpha = \gamma \cup \{(\lambda.x, \tau \text{ var})\}$. Since by assumption $\mu \sim_\gamma \nu$, LEMMA 1 gives us $\mu \cup \{(\lambda.x, \mu(e))\} \sim_\alpha \nu \cup \{(\lambda.x, \nu(e))\}$. Moreover, as the theorem is proven for well-typed commands that do not include functions, we can assert that $\mu'' \sim_\alpha \nu''$. Additionally, since we have $\text{dom}(\mu') = \text{dom}(\mu'') \setminus \{\lambda.x\}$ and $\mu'(e) = \mu''(\lambda.x)$, LEMMA 2 implies $\mu' \sim_\gamma \nu'$. Consequently, we conclude that $C' \sim_\gamma D'$. \square

The non-interference theorem says that when two configurations differ only on public variables, any code, well typed according to these rules, applied to these configurations yields configurations differing only in public variables. This ensures that modifying public variables, which could be manipulated by malicious actors, divulges no information about private data. There's a complete absence of leakage from private to public variables.

4 Implementation

Our implementation¹ begins by defining data types to represent the syntax of expressions, commands, and other language constructs. These data types capture the structure of the language and provide a foundation for type inference.

Besides, we represent a formal proof tree using an OCaml datatype. Each node in the tree corresponds to a step in the proof, and the tree structure captures the logical flow of the proof process.

The datatype defined in our implementation captures various derivations and subderivations within the proof tree, allowing us to construct a structured representation of the proof process. To understand more this dataset, consider the following code:

$$f ::= \text{fun } (\lambda.x) \{ \lambda.y := \lambda.x + 1; \text{return } \lambda.y \}; a := f(5)$$

Here is an example of a proof tree that demonstrates that this code is secure and therefore there cannot be any leak of private information into public variables:

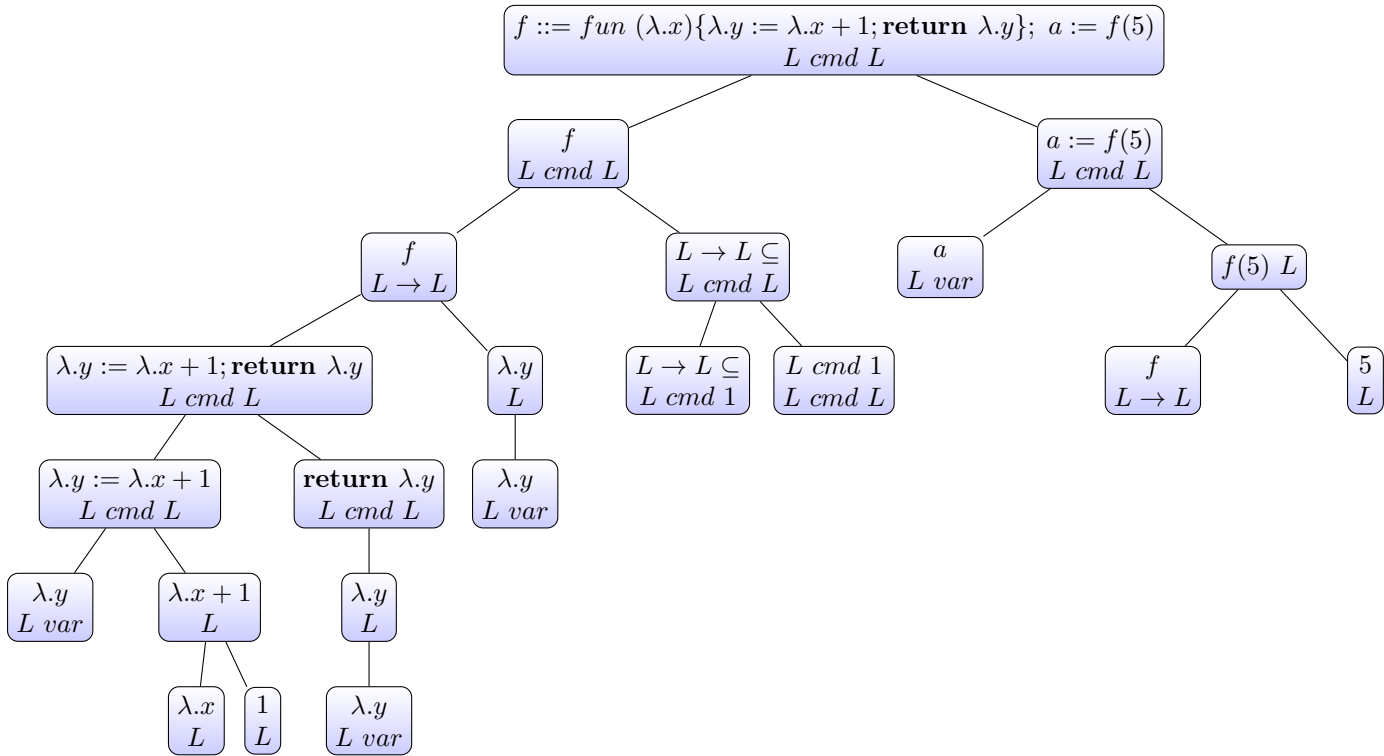


Figure 1: A formal proof tree

The proof tree datatype definition in our implementation is as follows:

```

type tree =
| Empty
| VarDeriv of string * phrase_type
| ConstDeriv of int
| BinopDeriv of phrase * phrase_type * tree * tree
| SkipDeriv of phrase_type
| AssignDeriv of string * phrase * phrase_type * tree * tree
| IfDeriv of phrase * phrase_type * tree * tree * tree
| WhileDeriv of phrase * phrase_type * tree * tree
| SeqDeriv of phrase * phrase_type * tree * tree
| ReturnDeriv of exp * phrase_type * tree
| FuncDefDeriv of string * string * command * phrase_type * tree

```

¹github.com/ymiled/secure-typing

```

| FuncNonVoidCallDeriv of string * exp * phrase_type * tree
| FuncVoidCallDeriv of string * string * command * phrase_type
(* For the subtyping rules : *)
| SubDeriv of phrase * phrase_type * phrase_type * tree * tree
| SubRules of phrase_type * phrase_type * tree * tree

```

This tree (1) is represented by the `example_proof_tree` (see below, 6).

We also implement a `type_checker` function (`check_type`) that verifies if a proof (represented in an OCaml formal proof tree as defined previously) is correct.

The ‘`check_type`’ function is the central piece of the implementation. It takes an inference tree, along with a context (`gamma`) which is a list describing the typing environment, and recursively traverses the tree to verify that types inferred at each node align with the language’s type system. The function also handles subtyping relationships and function calls, verifying that argument types match parameter types and that the returned type of a function matches the expected type.

When we test if the proof is correct with respect to the typing rules, i.e when we call

```
let verification = check_type example_proof_tree [("a", Var L)]
```

we get `true`. This means that the proof is correct (with respect to the typing rules), therefore by the non-interference theorem (1) the code `f ::= fun (λ.x){λ.y := λ.x + 1; return λ.y}; a := f(5)` is guaranteed to verify information flow security.

5 Conclusion

5.1 Summary

This paper explored a formal method to secure the information flow in a program, by using a specific type-based approach.

We extended existing typing rules and semantics for simple commands to handle the integration of functions, a very useful tool in a code. Through a theoretical framework, we were able to build a type system made to enforce security-centric typing regulations. We proved that the typing rules uphold the non-interference property, which asserts that when two configurations differ only on public variables, any code, well typed according to these rules, applied to these configurations yields configurations differing only in public variables. This ensures that modifying public variables, which could be manipulated by malicious actors, divulges no information about private data. There’s a complete absence of leakage from private to public variables.

We then developed a more practical approach by implementing and choosing certain representations to formally verify the well-typed nature of code. This involved the development of a type checker tasked with ensuring the correctness of a proof in our model.

5.2 Future directions

Our work presents potential development in the following areas :

- **Automatic type inference:** We are thinking about this approach, where types are automatically inferred, as a security information flow mechanism. A first step could be to develop a tool that help building a formal proof tree, by suggesting to the user, in a manner somewhat resembling the Coq language, how to proceed to prove that a phrase has a certain type.
- **Machine Learning for Type Checking:** One possible line of research is whether machine learning can help with the process of type checking. We may train machine learning models on secure programs in order to recognize patterns on formal proof trees and how to proceed.

References

- [1] FLEMMING NIELSON, HANNE RIIS NIELSON: Formal methods, an appetizer.
- [2] DAVID BASIN: Formal methods for security.
- [3] XAVIER LEROY: Le logiciel, entre l'esprit et la matière: Leçon inaugurale prononcée au Collège de France le jeudi 15 novembre 2018.
- [4] VINCENT SIMONET: Flow Caml in a Nutshell.
- [5] NEVIN HEINTZE, JON G. RIECKE: The Slam Calculus: Programming with Secrecy and Integrity.
- [6] GEOFFREY SMITH: A New Type System for Secure Information Flow.
- [7] DENNIS VOLPANO, GEOFFREY SMITH: A Type-Based Approach to Program Security.
- [8] GURVAN LE GUERNIC, THOMAS JENSEN: Monitoring Information Flow.

6 Appendix

The example proof tree (1):

```
let assign_tree =
  SubDeriv(
    Command (Assign("a", FuncCall("f", Int 5))),
    Cmd(L, L),
    Ncmd(L, 1),
    AssignDeriv(
      "a",
      Exp (FuncCall("f", Int 5)),
      Ncmd(L, 1),
      VarDeriv("a", T L),
      FuncNonVoidCallDeriv("f", Int 5, T L, ConstDeriv(5))
    ),
    SubRules(
      Ncmd(L, 1),
      Cmd(L, L),
      Empty,
      Empty
    )
  )

let function_command_tree =
  AssignDeriv(
    "y",
    Exp (Binop(Plus, Var "x", Int 1)),
    Ncmd(L, 1),
    VarDeriv("y", T L),
    BinopDeriv(
      Exp( Binop(Plus, Var "x", Int 1)),
      T L,
      VarDeriv("x", T L),
      ConstDeriv(1)
    )
  )

let func_tree =
  SubDeriv(
    Command (
      FuncDef(
        "f", "x", Seq( Assign("y", Binop(Plus, Var "x", Int 1)),
          (Return (Var "y")) ) ) ),
    Cmd(L, L),
    Func(L, L),
    FuncDefDeriv(
      "f",
      "x",
      Seq( Assign("y", Binop(Plus, Var "x", Int 1)), (Return (Var "y")) ),
      Func(L, L),
      function_command_tree
    ),
    SubRules(
      Func(L, L),
      Cmd(L, L),
      SubRules(
```

```

    Func(L, L),
    Ncmd(L, 1),
    Empty,
    Empty
  ),
  SubRules(
    Ncmd(L, 1),
    Cmd(L, L),
    Empty,
    Empty
  )
)
)
)

```

```

let example_proof_tree =
  SeqDeriv(
    Command (
      Seq (
        FuncDef(
          "f", "x", Seq( Assign("y", Binop(Plus, Var "x", Int 1)),
            (Return (Var "y")) )
          ),
          (Assign("a", FuncCall("f", Int 5)))
        )
      ),
    ),
    Cmd(L, L),
    func_tree,
    assign_tree
  )

```