

On Protection by Layout Randomization

Martín Abadi Gordon Plotkin

COLLOQUE ANNIVERSAIRE
EN L'HONNEUR DE
GÉRARD BERRY ET JEAN-JACQUES LÉVY

Outline

- 1 Introduction
- 2 Goals and Methodology
- 3 High Level
- 4 Low Level
- 5 Further Work

Topic

- Our main theme concerns formal proofs of security, both in computational models and in symbolic models.
- This talk is about low-level software protection, rather than cryptographic protocols.
- More could be done in this direction, perhaps learning from work on protocols.

Low-Level Attacks and Protection

- Many attack techniques:
 - Buffer overflows
 - Exception overwrites
 - return-to-libc
 - jump-to-libc
 - use-after-free attacks
- Often with knowledge and control of the heap
- Many defenses:
 - Stack canaries
 - Safe exception handling
 - NX (No eXecute) data
 - **Layout randomization**
- Useful mitigations
- But not necessarily perfect in a precise sense
- Nor all well understood

Executing Existing Code

- With NX defenses, attackers cannot simply inject data and then run it as code.
- But attackers can still run existing code:
 - an existing function, such as `system()`,
 - even code in the middle of a function,
 - even “accidental” code (e.g., starting half-way in a long x86 instruction),
 - even dead code.

Layout Randomization

- Attacks often depend on addresses.
⇒ *Let us randomize the addresses!*
 - Considered for data at least since the rise of large virtual address spaces (e.g., [Druschel & Peterson,1992] on fbufs).
 - Now present in Linux (PaX) and Windows (ASLR).

Goals

- Define a form of layout randomization.
- Study it precisely.
- Draw on ideas from the theory of programming languages and security, where:
 - there is much previous work on using equivalences for security properties (e.g., in information-flow control and in the study of protocols).
 - and there is work on preservation of security properties, perhaps given by such equivalences, by compilation.
- *Not* (yet) model and prove the security of actual systems.

Private and Public Variables

- We distinguish **private** variables (which should not be directly accessible to attackers) from **public** ones.
- Privacy gives secrecy and integrity properties:
 - If l is private, then $l := c$ should not reveal c .
 - If l is private, then an attacker should not be able to tamper with l in

```
 $\lambda f:\text{nat} \rightarrow \text{unit}.$   
   $l := c;$   
   $f(c);$   
   $\text{if } !l = c \text{ then } l' := c \text{ else } l' := c'$ 
```


Implementation

- l is mapped to a natural-number address $w(l)$, where w is a *layout* function.
- So, for example, $l := c$ compiles to $w(l) := c$

High- and Low-Level Threats

A high-level attacker may

- access public variables,
- call functions, pass arguments.

A low-level attacker may in addition access memory through natural-number addresses. There are two models:

- **Fatal Error Model** A fatal error is raised on an erroneous low-level memory address (as may happen in kernel mode).
- **Recoverable Error Model** A recoverable error is raised on an erroneous low-level memory address (as may happen in user mode).

Implementation (cont.)

- l is mapped to a natural-number address $w(l)$, where w is a *layout* function (w is fixed on public locations).
- So, for example, $l := c$ compiles to $w(l) := c$
- If an attacker can guess w , then it can break the expected security properties.
- Choosing w randomly may preserve those properties.
- We may expect strong guarantees for the fatal error model, but we may only expect guarantees parametrized on some count of memory accesses for the recoverable error model.

Approach

- Define high-level and low-level languages.
→ View the implementations as translations.
- View attackers as contexts.
→ Study whether low-level contexts correspond to high-level contexts.
- Phrase security properties as equivalences.
→ Study whether equivalences are preserved.
- In this talk we only consider the recoverable error model.

The Source Language

- Higher-order lambda calculus,
- with reading and writing operations on locations that hold natural numbers,
- and with standard base types and a type of locations.

Syntax: Types

- These are:

$$\sigma ::= \text{loc} \mid \text{nat} \mid \text{unit} \mid \sigma \times \tau \mid \sigma + \tau \mid \sigma \rightarrow \tau$$

- Commands are terms of type `unit`; they are evaluated only for their side-effects.
- We don't need a boolean type as we can take `bool` to be `unit + unit`.

Syntax: Terms

- Terms are ranged over by M, N, P, \dots

$M ::= x$	Variables
$ (l \in \text{Loc}) \mid ! \mid := \mid c$	Constants
$ (M, N) \mid \text{fst } M \mid \text{snd } M$	Products
$ \text{inl } M \mid \text{inr } M$	Sums
$ \text{cases } P \text{ inl } x:\sigma. M \text{ inr } y:\tau. N$	
$ \lambda x:\sigma. M \mid MN$	Functions
$ \text{letrec } f(x:\sigma) = M \text{ in } N$	Recursion

- Loc is a finite set of locations; it is divided into sets PubLoc and PriLoc of **public** and **private** locations .
- c ranges over arithmetic operations and predicates.
- Programs** are simply terms with no free variables.

Syntax: Typing

- Terms are given types $M:\sigma$ using typing axioms and rules.
- Some example typing axioms for constants:

$l:\text{loc} \ (l \in \text{Loc}) \quad !:\text{loc} \rightarrow \text{nat} \quad :=:\text{loc} \times \text{nat} \rightarrow \text{unit}$

- Some example typing rules

$$\frac{M:\sigma \quad N:\tau}{(M, N):\sigma \times \tau} \qquad \frac{M:\sigma \times \tau}{\text{fst } M:\sigma}$$

Contexts and Attackers

- A **context** is a program of type

$\sigma \rightarrow \text{bool}$

- A **public context** is one that cannot access private locations directly. That is, it contains no locations in `PriLoc`.
- For us:

attackers = public contexts

Full Abstraction

- Given a language L , two programs M and N are **observationally equivalent**

$$M \approx_L N$$

if they have the same **observable behavior** in all contexts.

- A translation $M \mapsto M^\downarrow$ from L to another language L' is **fully abstract** if it preserves and reflects observational equivalence:

$$M \approx_L N \quad \text{iff} \quad M^\downarrow \approx_{L'} N^\downarrow$$

Full abstraction (preview)

- We adapt these ideas to public contexts.
- We also weaken the low-level equivalences, with probabilities and complexity bounds.
- Modulo these changes, we do obtain full abstraction
 - The results are specific to certain languages.
 - They fail for some subsets and supersets.

Source Language: Semantics

- We define a **transition relation**

$$(s, M) \rightarrow (s', M')$$

by axioms and rules.

- Example axioms:

$$\begin{aligned} (s, !l) &\longrightarrow (s, n) & (s(l) = n) \\ (s, l := n) &\longrightarrow (s[l \mapsto n], \text{skip}) \end{aligned}$$

- Computations are then either **finite**:

$$(s_1, M_1) \rightarrow (s_2, M_2) \rightarrow \dots \rightarrow (s, V) \quad \equiv_{\text{def}} \quad (s_1, M_1) \Rightarrow (s, V)$$

where V is a value, e.g., l , 3 , $(4, \text{inl } 3)$, or $\lambda x:\text{nat}. x + x$.

- or **infinite**:

$$(s_1, M_1) \rightarrow (s_2, M_2) \rightarrow \dots \quad \equiv_{\text{def}} \quad (s_1, M_1) \uparrow$$

Source Language: Evaluation Function

- Evaluation functions give the observable behaviour of programs.
- **Eval** is the high-level evaluation function; it is defined by:

$$\text{Eval}(M, s) = \begin{cases} (s', V) & \text{if } (s, M) \Rightarrow (s', V) \\ \Omega & \text{if } (s, M) \uparrow \end{cases}$$

- So if termination is normal, Eval maps a program and an initial store to the final store and resulting value.

Source Language: Observational Equivalence

Two high-level programs are **publicly observationally equivalent** if no public context can distinguish them:

For M, N of type σ ,

$$M \approx_{h,p} N$$

iff for every public context $C: \sigma \rightarrow \text{bool}$ and store s ,

- (1) **either** both $\text{Eval}(CM, s)$ and $\text{Eval}(CN, s)$ diverge,
- (2) **or else** they both yield the same result value and two new publicly-equivalent stores (i.e., two stores that coincide on PubLoc).

Equivalences (cont.)

Secrecy and integrity properties can be phrased as public equivalences:

$$l := c \approx_{h,p} l := c'$$

and

$\lambda f:\text{nat} \rightarrow \text{unit}.$

$l := c;$

$f(c);$

$\text{if } !l = c \text{ then } l' := c \text{ else } l' := c'$

$\approx_{h,p}$

$\lambda f:\text{nat} \rightarrow \text{unit}.$

$l := c;$

$f(c);$

$l' := c$

The Target Model, Informally

- A **layout** w is a function of type $\text{Loc} \hookrightarrow \{0, \dots, c\}$ chosen uniformly at random.
- A **memory** m is a function of type $\{0, \dots, c\} \rightarrow \mathbb{N} + 1$
 - Memory may be accessed directly through natural-number addresses.
 - Some addresses may be unused.
- Accesses to unused addresses are **recoverable errors**

Target Language

- The types are:

$$\sigma ::= \text{loc} \mid \text{nat} \mid \text{unit} \mid \sigma \times \tau \mid \sigma + \tau \mid \sigma \rightarrow \tau$$

- The terms are as before, but the memory access constants are now typed differently:

$$\begin{aligned} l &:: \text{nat} \quad (l \in \text{Loc}) \\ ! &:: \text{nat} \rightarrow \text{nat} + \text{unit} \\ := &:: \text{nat} \times \text{nat} \rightarrow \text{com} + \text{unit} \end{aligned}$$

- The “+unit” represents the recoverable error case.
 We write `error` for `inr *`.

Exploring the Memory

Sample programs that explore memory:

```
cases !0 inl x:nat.x  
      inr y:unit.!1
```

```
cases !0 inl x:nat.x  
      inr y:unit.cases !1 inl x:nat.x  
                    inr y:unit.!2
```

→ Results for the recoverable error model require bounds on the number of accesses.

Our semantics therefore takes such accesses into account.

Semantics

- The semantics is given by axioms and rules for two transition relations, each relative to a layout:

$$w \models (m, M) \rightarrow (m', M') \quad w \models (m, M) \xrightarrow{a} (m', M')$$

- with the latter indicating a faulty memory address at a .
- Some semantics:

$$\begin{aligned} w \models (m, l) &\rightarrow (m, m(w(l))) \quad (l \in \text{Loc}) \\ w \models (m, !a) &\rightarrow (m, \text{inl } n) \quad (\text{if } a \in \{0, \dots, c\} \text{ and } m(a) = n) \\ w \models (m, !a) &\xrightarrow{a} (m, \text{error}) \quad (\text{if } a \notin \{0, \dots, c\} \text{ or } m(a) = *) \\ w \models (m, a := n) &\rightarrow (m[a \mapsto n], \text{inl skip}) \quad (\text{if } a \in \{0, \dots, c\} \text{ and } m(a) \neq *) \\ w \models (m, a := n) &\xrightarrow{a} (m, \text{error}) \quad (\text{if } a \notin \{0, \dots, c\} \text{ or } m(a) = *) \end{aligned}$$

(recalling that we write `error` for `inr *: nat + unit`)

Computations

- As before, computations are then either **finite**:

$$(m_1, M_1) \xrightarrow{(a_1)} \dots \xrightarrow{(a_{n-1})} (m_n, V) \equiv_{\text{def}} (m_1, M_1) \xrightarrow{A} (m, V)$$

where (a) indicates a possible a and A is the set of the actually occurring such a 's.

- or **infinite**:

$$(m_1, M_1) \xrightarrow{(a_1)} \dots \xrightarrow{(a_{n-1})} (m_n, V) \xrightarrow{(a_n)} \dots \equiv_{\text{def}} (m_1, M_1) \uparrow^A$$

Evaluation Function & Faulty Memory Accesses

- Eval_w^b is the low-level evaluation function; it has a layout parameter w , and a memory-access bound, b :

$$\text{Eval}_w^b(M, m) = \begin{cases} (m', V) & \text{if } w \models (m, M) \xrightarrow{A} (m', V) \\ & \text{and } |A \cap \{0, \dots, c\}| \leq b, \text{ for some } A \\ \Omega & \text{otherwise} \end{cases}$$

- When it terminates normally, and within bounds, it maps a program and initial memory to a final memory and value.
- It is also useful to record **the faulty memory accesses within bounds**, setting:

$$\text{Acc}_w(M, m) = A \cap \{0, \dots, c\}$$

where either $(m, M) \xrightarrow{A}$ – or $(m, M) \uparrow^A$.

Translations

- With each high-level program $M:\sigma$ we associate a low-level program $M^\downarrow:\sigma^\downarrow$.
Essentially this just replaces `loc` by `nat`.
- With each low-level context $C:\sigma \rightarrow \text{bool}$ we associate a high-level context $C^\uparrow:\sigma \rightarrow \text{bool}$.
Essentially this just replaces non-public memory accesses by recoverable errors.

Probability Bound

- δ_b is the probability that b distinct probes do not hit any private location of a layout chosen uniformly at random.
- We have:

$$\delta_b = \frac{|\mathbf{w}: \text{PriLoc} \hookrightarrow \{1, \dots, a - b\}|}{|\mathbf{w}: \text{PriLoc} \hookrightarrow \{1, \dots, a\}|}$$

where $a =_{\text{def}} (c + 1) - |\text{PubLoc}|$ is the amount of memory available for private locations (and assuming $b \leq a$).

- In case $|\text{PriLoc}| = 1$ this simplifies to:

$$\delta_b = 1 - b/a$$

Results for the Recoverable-Error Model: Theorem on Attackers

- Suppose $C:\sigma \rightarrow \text{bool}$ is a low-level attacker, i.e., a public low-level term (so σ is `loc`-free). We show that $C^\uparrow:\sigma \rightarrow \text{bool}$ is a high-level attacker of equal power, probabilistically.
- The theorem bounds the probability that the semantics of $C^\uparrow M$ and CM^\downarrow coincide, for a randomly chosen layout w .
- There is an assumption on the number b of erroneous accesses: without this bound, an attacker could explore all of memory.
- For small b , the high- and low-level semantics coincide the most, with probability close to 1 when c is sufficiently large.

Statement of Theorem

Theorem

Let $M:\sigma$ be a high-level term and let $C:\sigma \rightarrow \text{bool}$ be a low-level attacker. Then, for any store s , and $0 \leq b \leq c - |\text{Loc}|$, one of the following holds:

- 1 $\text{Prob}(|\text{Acc}_w(\text{CM}^\downarrow, s_w)| > b) \geq \delta_{b+1}$, or
- 2 $\text{Prob}(|\text{Acc}_w(\text{CM}^\downarrow, s_w)| \leq b \text{ and } \text{Eval}(C^\uparrow M, s)_w = \text{Eval}_w^b(\text{CM}^\downarrow, s_w)) \geq \delta_b$.

These alternatives are mutually exclusive if $\delta_{b+1} > 1/2$.

Remark: We would like to improve δ_{b+1} to δ_b here, and below, by considering *possibly* faulty memory accesses.

Low Level Partial Observational Equivalence

We set

$$M \approx_{l,p}^b N$$

iff for all low-level public contexts $C : \sigma \rightarrow \text{bool}$ and stores s :

- Either, for some $s' \upharpoonright \text{PubLoc} = s'' \upharpoonright \text{PubLoc}$ and V :

$$\text{Prob}(\text{Eval}_w^b(CM, s_w) = (s'_w, V)) \geq \delta_b$$

and

$$\text{Prob}(\text{Eval}_w^b(CN, s_w) = (s''_w, V)) \geq \delta_b$$

- or:

$$\text{Prob}(\text{Eval}_w^b(CM, s_w) = \Omega) \geq \delta_{b+1}$$

and

$$\text{Prob}(\text{Eval}_w^b(CN, s_w) = \Omega) \geq \delta_{b+1}$$

Theorem

Full Abstraction for Recoverable Error Model

Let $M, N: \sigma$ be high-level terms. Then, assuming that σ is loc-free , $0 \leq b \leq c - |\text{Loc}|$, and $\delta_{b+1} > 1/2$, we have:

$$M \approx_{h,p} N \quad \text{iff} \quad M^\downarrow \approx_{l,p}^b N^\downarrow$$

Further Work: Passing Locations

- Our results apply to passing functions that manipulate locations,
e.g., $(\lambda x:\text{unit}.!l, \lambda x:\text{nat}.l := x)$
- but not at all to passing locations,
e.g., two naked private locations l, l'

An Example (Poisoning)

An implementation could be “poisoned” with a natural number that does not represent a location, e.g.,:

$$M = \lambda x:\text{loc}.3$$
$$N = \lambda x:\text{loc}.\text{let } y \text{ be } !x \text{ in} \\ \text{cases } y \text{ in } \lambda x:\text{nat}.3 \text{ in } \lambda y:\text{unit}.5$$

are high- but not low-level publicly observationally equivalent.

An Example (Stateful Secrecy)

Here, in an implementation, an attacker could pass an f that stores its input, a secret location, for future use:

$$\begin{aligned}
 M &= \lambda f:\text{loc} \rightarrow \text{unit}. \\
 &\quad \text{if } !l_2 = 0 \text{ then } l_2 := 1; f(l_1) \text{ else } \Omega \\
 &\quad l_1 := 0 \\
 M' &= \lambda f:\text{loc} \rightarrow \text{unit}. \\
 &\quad \text{if } !l_2 = 0 \text{ then } l_2 := 1; f(l_1) \text{ else } \Omega \\
 &\quad l_1 := 1
 \end{aligned}$$

However at high level this is not possible, and as f can be called at most once, M and M' are publicly observationally equivalent.

Further Work (cont.)

- Our languages lack some features:
 - **code pointers**,
 - **dynamic allocation**,
 - ...
- Further work could discuss more attacks:
 - **jump-to-libc** attacks,
 - **use-after-free** attacks,
 - ...
- In another direction, further work could examine more features of implementations, e.g., **replication**.
- It would also be good to have methods to prove public equivalences $M \approx_{h,p} N$

*Félicitations à Gérard et
Jean-Jacques !*