



An Information Entropy Calculus for Stochastically Compiled Programs

Peter T. Breuer and Simon J. Pickin

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

November 29, 2019

An Information Entropy Calculus for Chaotically Compiled Programs

Peter T. Breuer¹ and Simon J. Pickin²

¹ Hecusys LLC, Atlanta, GA, USA. ptb@hecusys.com

² Universidad Complutense, Madrid, Spain. simon.pickin@fdi.ucm.es

Abstract. A calculus for the entropy in runtime traces resulting from stochastic program compilation is introduced here. It quantifies the variation that a ‘chaotic’ compiler aiming to vary the object code introduces into the program trace at run time, and the best strategy is characterised.

1 Introduction

This article describes stochastic compilation via a program calculus that quantifies the notion. To be clear from the outset, the stochastic element occurs not in the execution but in the compilation of a program. Our prototype HAVOC compiler (<http://sf.net/p/obfusc>) is for ANSI C [1] but the approach is generic. Our compiler’s particular constructions are guided by the principle:

Every machine code instruction that writes should introduce maximal possible entropy to the trace. (\tilde{h})

and, for that, a logic to reason about runtime entropy is required. The compiler’s aim in (\tilde{h}) is to vary object codes on each recompilation of the same source code so that runtime traces vary to the maximal extent possible. This kind of ‘chaotic’ compilation is first described in [2]. The program semantics is not conserved, but the difference from the nominal semantics is controlled. That is, the semantics obtained differs by a planned (different) ‘delta’ at every register and memory location before and after every instruction. That is called an *obfuscation scheme*. It makes the code and the runtime trace harder to read, so ‘obfuscation’ is an appropriate term for it. The scheme is a secret initially known only to the compiler and in principle shared only with the user/owner of the code. Reasons why it may be desirable to obfuscate will not be gone into here as introducing the technology is the focus, but they include intellectual property protection and other goals [3]. How it works may be illustrated via the following simple loop:

while $x < y + z + 1$ do { $x \leftarrow x + 2$; $x \leftarrow x + 3$; }

Imagine new program variables X, Y, Z and shift them by different deltas from the program variables x, y, z at different points in the code as shown below:

while $\underbrace{X+4}_x < \underbrace{Y+5}_y + \underbrace{Z+6}_z + 1$ do { $\underbrace{X+7}_x \leftarrow \underbrace{X+4}_x + 2$; $\underbrace{X+4}_x \leftarrow \underbrace{X+7}_x + 3$; }

The imagined relation $x = X + 4$ has to be the same at the end of the loop as at the beginning, but otherwise the choice is free. Simplifying:

while $X < Y + Z + 8$ do $\{X \leftarrow X - 1; X \leftarrow X + 6; \}$

(signed 2s complement comparison is translation-invariant). A user can see the first while loop execute and understand it as this while loop, working their imagination to make the substitutions. The lesson is that one program can be understood in many ways by an observer, who simply believes that the program variables are different by different deltas at different points in the program.

Conversely, a user intending the final while loop can instead obfuscate it to the first while loop, installing the given obfuscation scheme of deltas through the program, let it execute, and, discounting again the obfuscation, trust that it is the final while loop that has executed, which is really the one of interest.

The obfuscated codes (a) look approximately the same, having the same structure and differing only in program constants. Also (b) runtime traces ‘look the same’, with the same instructions in the same order reading and writing the same registers and memory locations (the end objective is to obfuscate at machine code level), while data varies from nominal by planned but arbitrary deltas, different at every point in the runtime trace and registers/memory, with the obvious provisos that:

deltas are equal across copy or skip, and where control paths meet. ($\underline{\text{h}}$)

That is necessary in order for computation to work properly. In particular, loops must have the same delta from nominal at either end, ready for a loop repeat.

This paper describes ‘correct by construction’ compilation for obfuscation as above following the principle ($\underline{\text{h}}$). For the compiled programs, at any m points in the trace not related as in ($\underline{\text{h}}$), it is shown that variations with $32m$ bits of entropy occur, supposing a 32-bit machine. That quantifies the obfuscation.

The compiler’s job is to vary the constants embedded in the machine code instructions so all feasible trace variations are exercised *equiprobably*. How it does that is summarised in Box 1: a new obfuscation scheme is generated at each recompilation. That is formally a set of vectors of *planned deltas from nominal values for the data per memory and register location, one such vector before and after each machine code instruction*. A declarative outline of the compile procedure is as follows: the compiler $\mathbb{C}[-]$ translates an expression e of type Expr that is to end up in register r at runtime to machine code mc of type MC and plans a 32-bit integer delta Δr (type Off) in r :

$$\begin{aligned} \mathbb{C}[-]^r &:: \text{Expr} \rightarrow (\text{MC}, \text{Off}) \\ \mathbb{C}[e]^r &= (mc, \Delta r) \end{aligned} \tag{1}$$

Let $s(r)$ be the value in register (or memory location) r in state s of the processor at runtime. The state is comprised by the values in registers and memory. Let $s(e)$ be the nominal (i.e., canonical, or at least standard, for some standard) evaluation of expression e . Running the code mc changes the state s to state s'

Box 1: An ‘obfuscating’ compiler does as follows:

- (A) *change only program constants, generating an arrangement of planned deltas from nominal values for instruction inputs and outputs (an obfuscation scheme);*
- (B) *leave runtime traces unchanged, apart from differences in the program constants (A) and runtime data;*
- (C) *equiprobably generate all arrangements satisfying (A), (B).*

that holds a value in r whose value differs by Δr from nominal. That is:

$$s \overset{mc}{\rightsquigarrow} s' \text{ where } s'(r) = s(e) + \Delta r \quad (2)$$

The compiler’s choices of deltas Δr before and after each machine code instruction are known to the user/owner of the code, but not the processor or the operating system. Since the user knows the obfuscation scheme, they can create meaningful inputs for the compiled program and interpret the outputs.

Section 2 introduces a concrete, modified OpenRISC (<http://openrisc.io>) machine code instruction set as a target for obfuscating compilation as described above. It can alternatively be understood as a set of assembly-level macros to be implemented by the instructions of another platform. The small-step instruction semantics in Table 1 concretizes what is described broadly by $s \overset{mc}{\rightsquigarrow} s'$ in (2).

The obfuscating compiler technology illustrated in the discussion of this section is described less abstractly via the translation of three canonical source code constructions in 3.1 and 3.2. A pre-/post-condition Hoare program logic [4] for the calculus of deltas that the compiler must use to keep track of its code variations is introduced in 3.3, and in 3.4 it is modified to an obfuscation calculus for stochastic compiler-induced entropy. A spin-off is that one is able to reason about statements such as the covariance $\text{cov}(x, x') = 0$ for a program variable x .

2 FxA Instructions

A ‘fused anything and add’ (FxA) [2] instruction set architecture (ISA) will be the specific compilation target here. The integer portion is shown in Table 1. The instruction set is adapted from OpenRISC ISA v1.1 (<http://openrisc.io/or1k.html>), which has about 200 instructions so FxA has about that many too. There are instructions for single and double precision integer operations, single and double floating point, and vector operations, all 32 bits long. Instructions access up to three 32 general purpose registers (GPRs), and one of those register operands may be replaced by a (‘immediate’) constant. A number of 32-bit ‘prefixes’ may precede a 32-bit instruction which has 16 bits of room for an on-board constant to provide room for 16 bits more of that constant and also any others that the instruction requires.

Table 1. FxA instruction set for encrypted work.

<i>op.</i>	<i>fields</i>	<i>mnem.</i>	<i>semantics</i>
add	$r_0 r_1 r_2 k$	add	$r_0 \leftarrow r_1 + r_2 + k$
sub	$r_0 r_1 r_2 k$	subtract	$r_0 \leftarrow r_1 - r_2 + k$
mul	$r_0 r_1 r_2 k_0 k_1 k_2$	multiply	$r_0 \leftarrow (r_1 - k_1) * (r_2 - k_2) + k_0$
div	$r_0 r_1 r_2 k_0 k_1 k_2$	divide	$r_0 \leftarrow (r_1 - k_1) \div (r_2 - k_2) + k_0$
...			
mov	$r_0 r_1$	move	$r_0 \leftarrow r_1$
beq	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 = r_2 + k$
bne	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 \neq r_2 + k$
blt	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 < r_2 + k$
bgt	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 > r_2 + k$
ble	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 \leq r_2 + k$
bge	$i r_1 r_2 k$	branch	if b then $pc \leftarrow pc + i$, $b \Leftrightarrow r_1 \geq r_2 + k$
...			
b	i	branch	$pc \leftarrow pc + i$
sw	$(k_0)r_0 r_1$	store	$\text{mem}[[r_0 + k_0]] \leftarrow r_1$
lw	$r_0 (k_1)r_1$	load	$r_0 \leftarrow \text{mem}[[r_1 + k_1]]$
jr	r	jump	$pc \leftarrow r$
jal	j	jump	$ra \leftarrow pc + 4$; $pc \leftarrow j$
j	j	jump	$pc \leftarrow j$
nop		no-op	
LEGEND			
r	- register index	k	- 32-bit integer
j	- prog. count	' \leftarrow '	- assignment
		i	- prog. incr.
		pc	- prog. count reg.
		ra	- return addr. reg.
		r	- register content

3 Obfuscating Compilation

An obfuscating compiler as outlined in Section 1 works with a database $D : \text{Loc} \rightarrow \text{Off}$ containing (here 32-bit) integer delta offsets Δl (type Off) for data, indexed per register or memory location l (type Loc). That is varied by the compiler as it makes its pass through the source code. The delta Δl defines by how much the runtime data is to differ from nominal in l at that point in the program control graph, and database D is an *obfuscation scheme* at that point.

The database $L : \text{Var} \rightarrow \text{Loc}$ that maps source code variables to registers and memory will not be treated here but taken as given (and annotated as a superscript on the compiler symbol).

Taking databases D (type DB) and L into account, the expression compiler $\mathbb{C}[e]^r$ described in Section 1 that places the result value in target register r in fact will be written $\mathbb{C}^L[D : e]^r$ with type signature:

$$\mathbb{C}^L[_ : _]^r : \text{DB} \times \text{Expr} \rightarrow \text{MC} \times \text{Off} \quad (3)$$

where MC stands for machine code, a sequence of FxA instructions mc . The compiler aims to vary the deltas Δl of type Off equiprobably over the type across recompilations. The following paragraphs explain how it is feasible.

3.1 Expressions

To translate $x + y$ where x, y are signed 32-bit integer source code variables, the compiler first emits machine code mc_1 as in (4a). At runtime that will put the

value of x in register $r_1 = Lx$ with offset delta Δx (a pair in $DB \times Expr$ is written $D : x$ for readability here):

$$(mc_1, \Delta x) = \mathbb{C}^L[D : x]^{r_1} \quad (4a)$$

$$s_0 \xrightarrow{mc_1} s_1 : s_1(r_1) = s_0(x) + \Delta x \quad (4b)$$

Small step semantics is from Table 1, with the value $s(r)$ in register r . To make that work formally, define the *nominal* value of a variable x in general as

$$s(x) =_{DF} s(Lx) - DLx \quad (5)$$

where $r = Lx$ is the location for the variable and Dr is the delta intended for that location at this point in the code. Then $\Delta r_1 = Dr_1$ and

$$s_1(x) = s_0(x) \quad (6)$$

has maintained the nominal value of variable x despite the numerical value in the register r_1 where it is located being different from that by Δx .

The compiler next emits machine code mc_2 . At runtime that will put the value of y in register $r_2 = Ly$ with offset delta Δy :

$$(mc_2, \Delta y) = \mathbb{C}^L[D : y]^{r_2} \quad (7a)$$

$$s_1 \xrightarrow{mc_2} s_2 : s_2(r_2) = s_1(y) + \Delta y \quad (7b)$$

As with x , that maintains the canonical value of y :

$$s_2(y) = s_1(y) \quad (8)$$

The compiler emits the FxA integer **add** instruction that at runtime adds the sum from r_1 and r_2 into r_0 , plus an increment k :

$$\mathbb{C}^L[D : x + y]^{r_0} = (mc_0, \Delta e) \quad (9a)$$

$$mc_0 = mc_1; mc_2; \mathbf{add} \ r_0 \ r_1 \ r_2 \ k$$

Setting the canonical value of a sum expression as $s(x + y) = s(x) + s(y)$, and choosing $k = \Delta e - \Delta x - \Delta y$, the expression gets the following value at runtime:

$$s_0 \xrightarrow{mc_0} s_2 : s_2(r_0) = s_0(x) + s_1(y) + \Delta e \quad (9b)$$

The register r_0 in which x is located will not be touched by the code mc_2 that evaluates y , because a compiler has always to be careful about that, so $s_1(x) = s_0(x)$ and the final value in r_0 is the nominal value of $x + y$ offset by delta Δe . The delta is independent of the two for x and y – the instruction constant k may be adjusted at will to suit.

3.2 Statements

Let Stat be the type of statements, then compiling a statement produces a new obfuscation scheme, as well as machine code:

$$\mathbb{C}^L[_ : _] : \text{DB} \times \text{Stat} \rightarrow \text{DB} \times \text{MC} \quad (10)$$

Consider an assignment $z=x+y$ of the previous expression $x+y$ to a source code variable z , which the location database L binds in register $r_z=Lz$. Let $x+y$ be called e here. The compiler emits code mc_0 that evaluates expression e in register $\mathbf{t0}$ with (randomly chosen) offset Δe as described in (9a) with $\mathbf{t0} = r_0$. A short-form **add** instruction with semantics $r_z \leftarrow \mathbf{t0} + k$ is emitted:

$$\mathbb{C}^L[D_0 : z=e] = D_1 : mc_0; \mathbf{add} \ r_z \ \mathbf{t0} \ k \quad (11a)$$

The compiler sets $k=\Delta r_z - \Delta e$ to choose delta Δz for z in $r_z = Lz$:

$$s_0 \xrightarrow{mc_0} s_2 \xrightarrow{\mathbf{add}} s_3 : s_3(r_z) = s_0(x) + s_1(y) + \Delta z \quad (11b)$$

The database of offset deltas is updated from $D_0 r_z$ to $D_1 r_z = \Delta z$ here, so that is

$$s_3(z) = s_0(x) + s_1(y) \quad (12)$$

Again, the final delta Δz may be freely chosen independent of the others by choosing the instruction constant k as required.

3.3 Offset Calculus

A classical pre-/post-condition calculus [4] captures more handily the compiler's changes to the obfuscation scheme of deltas in each register and memory location as it makes its pass through the source code.

Assignment. Generalising the $x+y$ above to expression e with intermediates in registers $\rho=\{r_0, \dots, r_n\}$, and result variable z stored in r_0 , the delta offsets before and after the assignment are generically:

$$\begin{aligned} & \{\Delta r_0 = Y_0, \dots, \Delta r_n = Y_n\} \\ & \quad \quad \quad z = e \\ & \{\Delta' r_0 = Z_0, \dots, \Delta' r_n = Z_n\} \end{aligned} \quad (13)$$

By the example (9b,11b), the $\Delta' r$, Δr are independently chosen as the compiler modifies (post-) vector Δ' to (pre-) vector Δ :

$$\{\Delta\} z = e \{\Delta'\} \quad (13a)$$

$$\text{where } \Delta \supseteq \Delta'|_{\bar{\rho}} \quad (13b)$$

The vectors Δ , Δ' are identical on the *complement* $\bar{\rho}$ of ρ .

Vectors. The Δ above are indexed by the entire type Loc of registers and memory locations but in practice only a small subset is needed for a program.

When pointers (memory addresses calculated dynamically) are involved we augment the type system of the source language so each is declared restricted to point into a particular global array as workspace:

```
int A[100]; ...; restrict A int *ptr;
```

That reduces the set of possible memory locations as indices of the Δ vectors (to the storage space for **A** in expressions involving ***ptr** here).

The predictable problem of not knowing at compile-time which entry in the array a pointer points to at run-time turns out not to be relevant: we are ultimately interested in the statistical variation of the content achieved by compiling stochastically, and it will ideally be uniformly the same for all array entries.

Conditionals. Source code conditionals are compiled to machine code branch instructions, but which branch is for true and which for false from boolean expression b can be deliberately varied by the compiler, which randomly chooses to generate either code for b or for $\neg b$ at each level of subexpression.

That compile procedure is detailed in [2], but it is not complicated. It has already been described: the 1/0 result b of each boolean subexpression is modified by a randomly chosen 1-bit delta δ to $b + \delta \pmod 2$ as just related for arithmetic expressions with the difference that the delta for booleans is 1-bit, not 32-bit.

The compiler tracks the nominal value, but it is not determinable by an observer whether the intended semantics has the nominal **then** or **else** branch as the target of the branch instruction jump, just as it is not determinable which **while** loop was meant by the code author in the introductory example given in Section 1. The same technique is used in classic ‘garbled circuits’ [5] technology for obfuscating hardware logic circuit design, with the difference that circuits cannot be recursive and boolean expressions can be.

The upshot is that the appropriate pre-/post- logic is a classic nondeterministic choice. Let ρ be the registers written in e . The deduction is:

$$\frac{\{\Delta_1\} s_1 \{\Delta'\} \quad \{\Delta_2\} s_2 \{\Delta'\}}{\{\Delta\} \text{ if } (e) s_1 \text{ else } s_2 \{\Delta'\}} \quad (14a)$$

$$\Delta \supseteq \Delta_1|_{\bar{\rho}} \cup \Delta_2|_{\bar{\rho}} \quad (14b)$$

and Δ , Δ_1 , Δ_2 are identical on $\bar{\rho}$, otherwise independent. The final offsets Δ' set by the compiler are equal at the end of both branches, as the code following it will be compiled supposing it receives the same scheme of offsets no matter which branch was executed.

Loops The compiler implements **do while** loops as body plus conditional branch back to the start. Let ρ be the registers written in e . The other registers

are equally offset at loop start and end, i.e., $\Delta_1|_{\bar{\rho}}=\Delta_2|_{\bar{\rho}}=\Delta'|_{\bar{\rho}}$ in (14a),(14b):

$$\frac{\{\Delta\} s \{\Delta'\}}{\{\Delta\} \text{ do } s \text{ while } e \{\Delta'\}} \quad (15a)$$

$$\Delta \supseteq \Delta'|_{\bar{\rho}} \quad (15b)$$

By this, the compiler is free to set offsets $\Delta|_{\rho}$ and $\Delta'|_{\rho}$ independently.

3.4 Obfuscation Calculus

Let f_r be the probability distribution of offset Δr from a nominal value v beneath the encryption in register r , as the compilation varies stochastically, so $\text{prob}(s(r) = v+d) = \text{prob}(\Delta r = d) = f_r(d)$, where s is the processor state.

Each Δr , $\Delta' r$ is a random variable with a probability distribution, giving the stochastic analogue below (16) of (13). Let variable x be stored in register or memory location r_x , and similarly for y , z ;

$$\begin{aligned} \{\Delta r_x = \mathcal{X}, \Delta r_y = \mathcal{Y}, \Delta r_z = \mathcal{Z}\} \\ z = x + y \\ \{\Delta' r_x = \mathcal{X}, \Delta' r_y = \mathcal{Y}, \Delta' r_z = \mathcal{Z}'\} \end{aligned} \quad (16)$$

In particular, Δr_z and $\Delta' r_z$ are independent random variables.

Let T be the runtime trace of a program. That is a sequence consisting of each instruction executed and the values it read and wrote. After an assignment the trace is longer by one: $T' = T \frown \langle z = e \rangle$.

The entropy $H(T)$ of the random variable T distributed as f_T is the expectation $\mathbb{E}[-\log_2 f_T]$. The increase in entropy from T to T' (it cannot decrease with longer T) is informally the number of bits of unpredictable information added.

We will need only these two facts from information theory:

Proposition 1 *The flat distribution $f_x=1/k$ constant is the one with maximal entropy $H(x)=\log_2 k$, on a signal x with k values.*

Proposition 2 *Adding a maximal entropy signal to any random variable on a n -bit space (2^n values) gives another maximal entropy, i.e., flat, distribution.*

Proposition 1 identifies the maximal entropy as n on an n -bit space, achieved when each of the 2^n values is equally probable. That is a completely disordered, or chaotic, signal. Proposition 2 uses the known fact (Shannon) that the entropy of the sum of two signals is at least as great as that of either. The surprising inference is that the characteristics of any distribution on a finite point space are obliterated completely, not partially, by adding a ‘chaotic’ signal to it, i.e., one with flat, uniform distribution.

Below, the logic is worked through for this stochastic view of compilation for the three constructs already seen: assignment, conditionals, and **while** loops.

Assignment. As in (13a), for pre-/post-condition:

$$\{\Delta\} z = e \{\Delta'\} \quad (17a)$$

but the Δ, Δ' are vectors of offsets $\Delta r, \Delta' r$ that are random variables as in (16).

Let $\rho = \{r_0, \dots, r_n\}$ be the registers written in e or in writing to z . For $r \notin \rho$, $\Delta' r = \Delta r$, because they are equal values by (13b), so the condition is still $\Delta|_{\tilde{\rho}} = \Delta'|_{\tilde{\rho}}$ here. I.e.:

$$\Delta \supseteq \Delta'|_{\tilde{\rho}} \quad (17b)$$

but there is more to come, because we suppose the compiler follows the principle $(\tilde{\mathbb{h}})$ and that means each new random variable is independent with maximal entropy. Each represents the compiler's free choice of constant like k of (9a,11a) in 'an arithmetic instruction that writes' (**add** and **addi** respectively in those).

Let the trace entropy up to the assignment be $H(T) = h$. Writing to r_z with offset delta that is a new independent r.v. U increases the trace entropy to $H(T') = h + H(U)$. The offset is 32-bit, chosen with flat distribution by the compiler, as per $(\tilde{\mathbb{h}})$, so $H(U) = 32$. There are $n+1$ registers r_0, \dots, r_n that are written independently, including the one holding target variable z , so entropy increases by $32(n+1)$ bits:

$$\{H(T) = h\} z = e \{H(T') = h + 32(n+1)\} \quad (17c)$$

In more general form, writing $\Phi[A/B]$ for substitution of A for B in predicate Φ :

$$\{\Phi[H(T) + 32(n+1)/H(T)]\} z = e \{\Phi'[H(T')]\}$$

Put $H(T) = h + 32(n+1)$ for Φ to obtain (17c) again.

That is the common case. But where the instruction that writes has already appeared once earlier in the trace, the offset delta it introduces is already known, and the increment in trace entropy is zero this time:

$$\{H(T) = h\} z = e \{H(T') = h\} \quad (17c0)$$

Alternatively:

$$\{\Phi[H(T)]\} z = e \{\Phi'[H(T')]\}$$

Remark 1. Whatever the expression e , provided it contains some arithmetic (even $z + 0$ will do), z' is independent of z is implied, as stated in Section 1.

Here ' z ' is understood to refer to the exact value in the register r_z where z is located. By (17c), If $H(T) = h$ before the assignment, then it is $H(T') = h + 32(n+1)$ after the assignment, where $n+1$ is the number of 'arithmetic instructions that write' that appear in the machine code, including one that last writes to r_z . That can only happen if each such instruction introduces 32 bits of entropy to the trace because that is maximal per instruction, and therefore the last write to r_z introduces 32 bits of entropy. That means it adds an independent

uniformly distributed 32-bit delta δ . Then $\text{prob}(r_{\mathbf{z}} = i \wedge r'_{\mathbf{z}} = j) = \text{prob}(r_{\mathbf{z}} = i \wedge r'_{\mathbf{z}} - r_{\mathbf{z}} = j - i) = \text{prob}(r_{\mathbf{z}} = i)\text{prob}(r'_{\mathbf{z}} - r_{\mathbf{z}} = j - i | r_{\mathbf{z}} = i) = \text{prob}(r_{\mathbf{z}} = i)\text{prob}(\delta = j - i | r_{\mathbf{z}} = i) = \text{prob}(r_{\mathbf{z}} = i)/2^{32} = \text{prob}(r_{\mathbf{z}} = i)\text{prob}(r'_{\mathbf{z}} = j)$ since $r'_{\mathbf{z}}$ is distributed uniformly too (as the sum with a maximal entropy difference signal δ). So $r_{\mathbf{z}}$ and $r'_{\mathbf{z}}$ are independent. and their covariance is 0, i.e., $\text{cov}(\mathbf{z}, \mathbf{z}') = 0$.

Conditionals. As in (13b),(14b) but with random variables:

$$\frac{\{\Delta_1\} s_1 \{\Delta'\} \quad \{\Delta_2\} s_2 \{\Delta'\}}{\{\Delta\} \text{ if } (b) s_1 \text{ else } s_2 \{\Delta'\}} \quad (18a)$$

$$\Delta \supseteq \Delta_1|_{\bar{\rho}} \cup \Delta_2|_{\bar{\rho}} \quad (18b)$$

The $\Delta r = \Delta_1 r = \Delta_2 r$ for $r \notin \rho$, because they are equal values according to (14b). The entropy added to the trace T is from the trace of b , say $32n$ bits of entropy from n writes to n registers, plus the entropy from the trace through a branch:

$$\frac{\{\text{H}(T')=h+32n\} s_1 \{\Theta\} \quad \{\text{H}(T')=h+32n\} s_2 \{\Theta\}}{\{\text{H}(T) = h\} \text{ if } (b) s_1 \text{ else } s_2 \{\Theta\}} \quad (18c)$$

or, in more generic form:

$$\frac{\{\Phi[\text{H}(T')]\} s_1 \{\Theta\} \quad \{\Phi[\text{H}(T')]\} s_2 \{\Theta\}}{\{\Phi[\text{H}(T) + 32n/\text{H}(T)]\} \text{ if } (b) s_1 \text{ else } s_2 \{\Theta\}}$$

Put $\text{H}(T) = h + 32n$ for Φ in the above to get (18c).

To make that deduction valid, the compiler must even up the number of arithmetic writes between the two branches so the entropy increase is the same. It can do it, because, even for loops, the entropy increase is finite and bounded (see below).

A second time the conditional appears in the trace, if it branches the same way again then it contributes zero entropy as the offset deltas are known:

$$\{\text{H}(T) = h\} \text{ if } (b) s_1 \text{ else } s_2 \{\text{H}(T') = h\} \quad (18c0)$$

If it branches a different way from the first time, the branch (but not the test) contributes entropy, as the offsets in that branch are yet unknown. But the, say m , instructions that align final offsets are constrained in (14b) to agree with the offsets in the other branch, which are already known. So those m do not count:

$$\frac{\{\text{H}(T)=h\} s_1 \{\Theta\} \quad \{\text{H}(T)=h\} s_2 \{\Theta\}}{\{\text{H}(T)+32m = h\} \text{ if } (b) s_1 \text{ else } s_2 \{\Theta\}} \quad (18c1)$$

Those final m instructions that synchronise the offsets with the other branch get a special name:

Definition 1 *An instruction emitted to adjust the final offset to a common value with the other branch is a trailer instruction. Each is last to write to a register in the branch.*

Loops Let $\rho = \{r_1, \dots, r_n\}$ be the registers written in b . Then, per (15a), (15b), but with random variables:

$$\frac{\{\Delta\} s \{\Delta'\}}{\{\Delta\} \text{ do } s \text{ while } (b) \{\Delta'\}} \quad (19a)$$

$$\Delta \supseteq \Delta'|_{\rho} \quad (19b)$$

That means $\Delta r = \Delta' r$ for $r \notin \rho$. Those distributions are equal because the values are equal for $r \notin \rho$, by (17b).

A trace over the loop is always the same length between recompilation and recompilation, because the compiler varies data values, not semantics. Say the loop repeats $N \geq 1$ times for a particular set of input values. Then it could be unrolled to N instances of the loop body and N instances of the loop test. The variation in the trace is only that of (a) the test repeated once, because the same offsets are applied to the n registers that are written in b at each repeat, plus (b) that of the body repeated once, for the same reason. The entropy calculation is (a) plus (b), no matter what N is:

$$\frac{\{H(T) + 32m = h\} s \{H(T') = h\}}{\{H(T) + 32(n+m) = h\} \text{ do } s \text{ while } b \{H(T') = h\}} \quad (19c)$$

Put $H(T) + 32m = h$ for Φ in the following generalisation to get (19c):

$$\frac{\{\Phi'[H(T')]\} s \{\Theta\}}{\{\Phi[H(T) + 32n/H(T)]\} \text{ do } s \text{ while } b \{\Theta\}}$$

The abstraction here is that a **do while** may lengthen the trace arbitrarily like a loop but it adds entropy to it like a conditional.

On a second time through the loop, zero entropy is added, because the offsets are the same as the last time:

$$\{H(T) = h\} \text{ do } s \text{ while } b \{H(T) = h\} \quad (19c0)$$

The (red) equations are an obfuscation calculus for trace entropy when compilation follows the principle ($\tilde{\mathbb{h}}$). In summary, counting up via the rules above:

Lemma 1 *The entropy of a trace is $32(n+i)$ bits; n is the number in it of distinct arithmetic instructions that write (a pair of trailer instructions count as the same) and i is the number of inputs.*

‘Inputs’ are identified with those instructions that read first time a location that has not yet been written in the trace.

A successfully obfuscating compiler’s constructions must recruit every emitted arithmetic instruction that writes to the task of freely varying the offsets in data beneath the encryption in register and memory locations. Otherwise it is not doing as much as it could to contribute to variability in the trace. The sole restriction on the compiler is that two final writes in different control paths must

set up the same offsets, and that is for correct program working. Conditionals, loops and gotos would go wrong otherwise. The following characterises the best way of compiling to produce this kind of obfuscation:

Proposition 3 *The entropy in the run-time traces induced by a compiler following the principle $(\tilde{\text{h}})$ is maximal among compositional strategies for varying the constant parameters in the compiled machine code while it still works correctly.*

The reasoning is that if the compiler works compositionally then it does not know the context in which its constructions are used, so it must suppose that data that is written will later be read, and so must arrange for synchronisation between all final offset deltas in different branches of conditionals, even for locations that in fact are never read afterwards. If the data is never read, synchronisation could be done without and total entropy would be increased when both branches are traversed in the trace. But for working compositionally, a compiler could put more entropy into the trace, but it would be for data that serves no purpose because it is never read.

The only other way to put more entropy in is to vary individual instructions more, but that is impossible for a compiler that already implements $(\tilde{\text{h}})$.

The proposition implies that on a 32-bit platform a full 32 bits of entropy per datum are provided by the compiler that follows $(\tilde{\text{h}})$:

Corollary 1 *The probability across different compilations that any particular 32-bit value x is in a given register or memory location at any given point in the trace at runtime is uniformly $1/2^{32}$.*

But what of two or more data values at different points in the trace? That depends on how they are connected computationally. If they are input and output of a copy instruction, they will be correlated.

Definition 2 *Two data values in the trace are (obfuscation) dependent if they are from the same register or memory location at the same point, are input and output of a copy instruction, or are from the same register or location at a join of two control paths after the last write to it in each and before the next write.*

If data is taken at two (m) independent points, the variation is maximal:

Theorem 1 *The probability across different compilations that any m particular 32-bit values have values x_i in given register or memory locations at given points in the program at runtime, provided they are pairwise independent, is $1/2^{32m}$.*

Each dependent pair reduces the entropy by 32 bits. The theorem asserts that with $(\tilde{\text{h}})$ as a guide, as much trace data as one cares to observe is maximally unpredictable across recompilations, as far as is computationally possible.

Remark 2. The theorem speaks to an argument about obfuscation in plain sight. It is impossible for any observer to know what the variables X, Y, Z are in the **while** loop given in Section 1, because they exist only in the mind of the code author. The problem is that the observer has only to duplicate the form of the

loop and run it to get a sequence of values x_0, x_1, x_2 , etc. in x and y and z that are duplicated at some unknown offset as X_0, X_1, X_2 , etc. by X, Y, Z . Thus $X_1 - X_0, X_2 - X_1$, etc. are known to the observer, being equal to respectively $x_1 - x_0, x_2 - x_1$, etc. The theorem says that if the code author is to avoid the observer being able to know even that much, then code with loops is out.

Suppose an observer claims to have a polynomial time (in the number of bits n in a word on the platform, hitherto supposed to be 32) method of working out what the code author means to have in their variable X at some identified point in the trace of program P . The point of interest to the observer may even move (polynomially) with n , being specified, say as ‘the point of last change in the first n^3 steps’. The observer knows program P and may even have suggested it. The observer also can see the compiled maximal entropy code $\mathbb{C}[P]$, and see it running and probe it deeply by running it himself.

The code author then readies a sequence of maximal entropy compilations $\mathbb{C}[P_n]$ of P , with the n th being for a n -bit platform as target, and with P having been partially or completely unrolled as P_n with no loops in at least the first e^n (i.e., super-polynomially many) machine code instructions. If the program predictably ends before then, just unroll completely. The observer is invited to apply their method and predict what is meant by the values at their chosen points in the trace(s) of these programs, which differ only in consequence of the number n of bits in a word on the platform .

The theorem implies the observer’s method cannot exist. There are ‘no loops’ (i.e., no dependencies, per the wording of the theorem) in the part of the trace the observer has time to examine. The theorem says the compiler will have arbitrarily and independently varied what is meant by the values throughout that length of the trace by varying the deltas from the nominal value independently across each instruction in turn. The observer, seeing a 1, cannot tell if 2 was meant.

The credibility of the argument is supported by the trivial case in which the program unrolls completely. Then it is equivalent to a logic circuit in hardware. It is known from the theory of Yao’s garbled circuits [5] that the intended values cannot be deciphered without the garbling scheme, which equates to an obfuscation scheme of deltas here via ‘+1 mod 2’ being boolean negation on a 1-bit logic value, while ‘+0 mod 2’ leaves the logic value unchanged.

4 Implementation

Our own prototype compiler <http://anonymised.url> following 3.4 is for ANSI C [1], where pointers and arrays present particular difficulties. Currently, the compiler has near total coverage of ANSI C and GNU C extensions, including statements-as-expressions and expressions-as-statements, gotos, arrays, pointers, structs, unions, floating point, double integer and floating point data. Pointers are obligatorily declared via ANSI **restrict** to point into arrays. It is missing **longjmp** and efficient strings (**char** and **short** are same as **int**), and global data shared across code units (a linker issue). The largest C source compiled (correctly) so far is 22,000 lines for the IEEE floating point test suite at [http:](http://)

Table 2. Trace for Ackermann(3,1), result 13.

PC	instruction	trace update
...		
35	add t0 a0 zer -86921031	t0 = -86921028
36	add t1 zer zer -327157853	t1 = -327157853
37	beq t0 t1 2 240236822	
38	add t0 zer zer -1242455113	t0 = -1242455113
39	b 1	
41	add t1 zer zer -1902505258	t1 = -1902505258
42	xor t0 t0 t1 -1734761313	1242455113 1902505258 t0 = -17347613130
43	beq t0 zer 9 -1734761313	
53	add sp sp zer 800875856	sp = 1687471183
54	add t0 a1 zer -915514235	t0 = -915514234
55	add t1 zer zer -1175411995	t1 = -1175411995
56	beq t0 t1 2 259897760	
57	add t0 zer zer 11161509	t0 = 11161509
...		
143	add v0 t0 zer 42611675	v0 = 13
...		
147	jr ra	# (return 13 in v0)

Legend: (registers) a0 = function argument; sp = stack pointer; t0, t1 = temporaries; v0 = return value; zer = null placeholder.

[//jhauser.us/arithmetic/TestFloat.html](http://jhauser.us/arithmetic/TestFloat.html). A trace³ of the Ackermann function⁴ [6] is shown in Table 2.

5 Conclusion

A formal obfuscation calculus for programs is set out that quantifies the entropy in the data beneath the encryption in a runtime trace, where compilation is stochastic (and execution is not). The way to maximise the entropy is to follow the principle for compiler constructions that every arithmetic instruction that writes should be varied maximally across recompilations.

References

1. ISO/IEC, “Programming languages – C,” Int. Org. for Standardization, 9899:201x Tech. Rep. n1570, Aug. 2011, JTC 1, SC 22, WG 14.
2. P. Breuer, J. Bowen, E. Palomar, and Z. Liu, “On obfuscating compilation for encrypted computing,” in *Proc. 14th Int. Conf. Sec. Crypto. (SECRYPT’17)*, P. Samarati *et al.*, Eds. SCITEPRESS, 2017, pp. 247–54.
3. B. Barak, “Hopes, fears, and software obfuscation,” *Commun. ACM*, vol. 59, no. 3, pp. 88–96, Feb. 2016.
4. C. A. R. Hoare, “An axiomatic basis for computer programming,” *Commun. ACM*, vol. 12, no. 10, pp. 576–80, 1969.
5. A. C.-C. Yao, “How to generate and exchange secrets,” in *27th Ann. Symp. Found. Comp. Sci.* IEEE, 1986, pp. 162–167.
6. Y. Sundblad, “The Ackermann function: a theoretical, computational, and formula manipulative study,” *BIT Num. Math.*, vol. 11, no. 1, pp. 107–19, 1971.

³ For readability here, the final offset delta for register **v0** is set to zero.

⁴ Ackermann C code: **int** A(**int** m,**int** n) { **if** (m == 0) **return** n+1; **if** (n == 0) **return** A(m-1, 1); **return** A(m-1, A(m, n-1)); }.