**UNIVERSITÄT**
**D U I S B U R G**
**E S S E N**

*Open*-*Minded*

# *Side-Channels in Cryptographic Software, the Haskell* `case`

attacks, semantics, CT criterion, and compiler construction

Marcel Fourné ■ 2020-01-14

- studied both (pure) Informatics and IT-Security

- self-taught Haskeller since 2005

- free software Haskell ECC

- https://hackage.haskell.org/package/eccrypto

  - NIST Prime Curves

  - Ed25519

  - utility code, easy to modify

  - predecessor library has been used for Ripple prototype, alternative Bitcoin etc.

- some work in the Debian Haskell group

- some work in the Haskell cryptography community

# Cryptography, abstract

1  op :: SecKey —> Data —> Result

- secret (key)
  - size of secret may be public knowledge
  - content of the secret must remain unknown to third parties
- data (not secret)
- some operation, a function
  - takes time and computational resources to compute
    - may be observable, depending on the attacker (model)
    - may have other observable(!) side effects
  - total functions are most commonly used
- result (maybe public)

# Side-Channels

Basically:

*Any* observation on a computation which allows

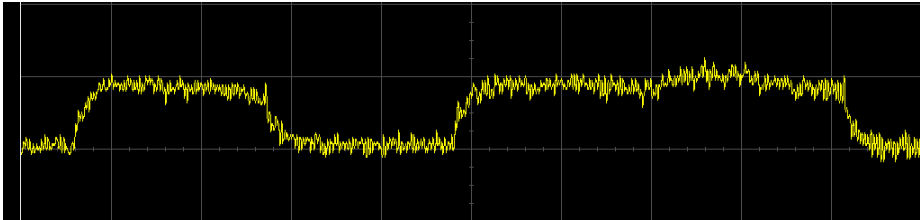inferences on the *content* of the secret key.

- Observations must be quantifiable.

- Inferences must be computable, but may be probabilistic.

    - partial bits

- 1 bit may be enough

- observations may be repeated

- may be preceded by an active attack

- unexpected results, software/hardware induced errors
  - wrong or specifically crafted results, error messages
  - correct results after a different time

- timing information (Kocher, 1996; Lipton&Naughton, 1993)
  - absolute values
  - relative differences

- power consumption (Kocher, 1998)
  - power consumption patterns (Simple Power Analysis)
  - power consumption pattern differences (Differential Power Analysis)

- memory access patterns and cache state (Percival, 2005; Bernstein, 2005)

- micro-architectural state (numerous since 2017, mainstream media famous since January 2018)

```
1   if(secretKeyBits[n] == 1) {
2       a * square(b)
3   }
4   else {
5       a * b
6   }
```

```
1  if(secretKeyBits[n] == 1) {
2      a * square(b)
3  }
4  else {
5      a * b
6  }
```

```
1  if(secretKeyBits[n] == 1) {
2      a * square(b)
3  }
4  else {
5      a * square(c)
6  }
```

```
1  if(secretKeyBits[n] == 1) {
2      a * square(b)
3  }
4  else {
5      a * square(c)
6  }
```

```
1   if(secretKeyBits[n] == 1) {
2      a * square(1)
3   }
4   else {
5      a * square(0)
6   }
```

```
1   if(secretKeyBits[n] == 1) {
2       a * square(1)
3   }
4   else {
5       a * square(0)
6   }
```

1    a + multiplicationTable [secretKeyBytes[n]]

1    a +  multiplicationTable [secretKeyBytes[n]]

# CONSTANT TIME Criterion

We need to prove two things:

- branch free in its secrets (PROGRAM COUNTER Model)

- no secret value dependent address indices

  ... but that is hard in a Turing Complete Programming Language

- *suggests* constant run-time behaviour, but:

  - garbage collection (no secrets)

  - operating system interaction and other nondeterministic behaviour

- may be overly prohibitive (see: non-CT variants of Montgomery Multiplication), but:

  - formulae can often be changed to be total and branch free

  - algebraic style line code instead of lookup tables or masked lookups

- may miss hardware side effects not dependent on program structure

- CT only needed for code handling secret values

- higher level code not affected (if type system allows no memory accesses across boundaries)

    - composability, linking, address spaces, see C and its common exploits

- operating system interactions

- hardware effects

    - known-bad Assembly instructions

    - value dependent latency multipliers on certain archs

    - microcode implementation replacements

Results:

- anything observable which is not yet usable as a side channel may still become one

- How to prove CT program behaviour? Which program behaviour is still permissible?

- correct by careful programming in Assembly

- manual analysis

- model in proof language, check with its type system

- generate implementation code from proof language

- check in implementation language type system

- correct by careful programming in Assembly

  - no compiler, implementation effort for each platform

- manual analysis

  - error-prone

- model in proof language, check with its type system

  - equivalence of model to implementation, assumptions

- generate implementation code from proof language

  - different language semantics and compilers

- check in implementation language type system

  - often not as expressive as dedicated proof languages

Compilers may change optimisations, but their type system guarantees should hold.

# *Haskell*

- non-strict evaluation semantics (graph reduction)

- garbage collection, strictness analysis, HM-style type inference, expressive code

- pure functions, explicit side-effects, easy parallelisation, cross-module inlining

- efficient code can be generated with Instructions Per Cycle $> 2$; benchmarks rival optimised C

- used in industry (Facebooks anti-spam, Bluespec System Verilog, seL4. . . )

- influences other languages (STM, list comprehensions, monads, QuickCheck)

- active type system research community, committed to correctness, very friendly

- non-strict evaluation semantics (graph reduction)

- garbage collection, strictness analysis, HM-style type inference, expressive code

- pure functions, explicit side-effects, easy parallelisation, cross-module inlining

- efficient code can be generated with Instructions Per Cycle $> 2$; benchmarks rival optimised C

- used in industry (Facebooks anti-spam, Bluespec System Verilog, seL4...)

- influences other languages (STM, list comprehensions, monads, QuickCheck)

- active type system research community, committed to correctness, very friendly

But:

- harder to reason about resource usage (garbage collection and lazy evaluation)

- reputation as hard to learn/too research-centric language

```
1    f 0 x _ = x
2    f n x y = f (n−1) y x
3
4    main = do
5       x <− readLn
6       let y = f x (ackermann 4 2) (ackermann 0 1)
7       print y
```

let: just like stating some unordered lemmas in mathematics before using them in a formula

```
1    g x = let i = m
2              k = x + 2
3              m = k + 2
4          in i * k * m
```

case: need to evaluate its condition before choice of branch

```
1    h x = case x of
2              0 −> 0
3              1 −> x + 1
4              _ −> x + 2
```

- Pedersen, Askarov, "From trash to treasure: timing-sensitive garbage collection" S&P'17 paper
  - allocates arrays differently based on secret values
  - requires non CT code

- in CT code memory may not be accessed dependent on secret value content, same for allocations

- CT code uses secrets with non-differing memory access patterns

- CT code is not vulnerable against this attack

*Parsing*, *Renaming*, *Typechecking*, and *Desugaring*, which produces `core`

⇓

*Simplification* works on `core` and it is where most optimisations happen

⇓

`core` is a minimal Haskell with explicit types; `case` for evaluation and branching, `let` for allocation

⇓

*STG* has thunks for laziness in `let`; three local optimisations after that

⇓

*Cmm* adds an explicit stack

⇓

*RTS.a* + the generated *Assembly*, which we analysed to check CT preservation

- total functional programming implies same behaviour between strict and lazy evaluation

- intuition: if function $f$ is called, then the result of $f$ is produced by computation

- without branches (etc.) in $f$: if $f$ is called, then CT behaviour happens in $f$

- necessary condition: no use (inspection) of secret values outside CT context

    - idea: have a verifiable subset of the program code and contain secret keys to this subset

    - implies control over copying during optimisations, which may be overly optimistic

- optimisations must be contained

    - Continuation Passing Style makes control flow explicit, but confuses IDA

- low-level calling conventions make more problems for manual analysis than evaluation order

As an example, why lazy evaluation can make non-CT code easier to exploit:

```
1    pmul :: EC -> Point -> Integer -> Point
2    pmul curve@(ECi l _ p _) b k  =
3        let ex p1 p2 i
4            | i < 0 = p1
5            | condBit k i == 0 = ex (pdouble curve p1) (padd curve p1 p2) (i - 1)
6            | otherwise     = ex (padd curve p1 p2) (pdouble curve p2) (i - 1)
7        in ex b (pdouble curve b) (log2len k - 2)
```

As an example, why lazy evaluation can make non-CT code easier to exploit:

```
1   pmul :: EC -> Point -> Integer -> Point
2   pmul curve@(ECi l _ p _) b k  =
3       let ex p1 p2 i
4             | i < 0 = p1
5             | condBit k i  == 0 = ex (pdouble curve p1) (padd curve p1 p2) (i − 1)
6             | otherwise      = ex (padd curve p1 p2) (pdouble curve p2) (i − 1)
7       in ex b (pdouble curve b) (log2len k − 2)
```

# manual proof, absence of required code gadgets

Also, how to find this at the assembly level:

```
1   pmul :: EC −> Point −> Integer −> Point
2   pmul curve@(ECi l _ p _) b k  =
3     let ex p1 p2 i
4         | i < 0 = p1
5         | condBit k i == 0 = ex (pdouble curve p1) (padd curve p1 p2) (i − 1)
6         | otherwise  = ex (padd curve p1 p2) (pdouble curve p2) (i − 1)
7     in ex b (pdouble curve b) (log2len k − 2)
```

Compiles to this comparison for the branch criterion:

```
1   mov r8, [rbx+rdi∗8+10h]
2   test r8, r8
3   jnz loc_3780
```

The address in line 1 contains a key content derived value, so our analysis flagged this branch non-CT.

- know your code generator and run-time to prevent unintentional miscompilation of security mechanisms

  - type check at which level? pre desugaring? post desugaring?
  - after `core`: some graph optimisations on `STG` (Shared Term Graph/Spineless, Tagless G-machine)
  - GHC-specific: Tables-Next-To-Code, `ghc-asm.lprl`, runs after Assembly code generator
  - modules with integrated C code, change of memory model

- know your code generator and run-time to prevent unintentional miscompilation of security mechanisms

  - type check at which level? pre desugaring? post desugaring?
  - after `core`: some graph optimisations on `STG` (Shared Term Graph/Spineless, Tagless G-machine)
  - GHC-specific: Tables-Next-To-Code, `ghc-asm.lprl`, runs after Assembly code generator
  - modules with integrated C code, change of memory model

- "When Constant-Time Source Yields Variable Time Binary: Exploiting Curve25519-donna Built with MSVC 2015"

  - verified C code executed with dependencies on unverified run-time libraries yielded vulnerable code

- branch-freeness via Abstract Syntax Tree, type annotated functions

- find type annotation at low-level Intermediate Language (`core`)

- build AST of function using Template Haskell and find problematic constructs, check at compile-time

- separation of concerns if types are enforced to be created only in some set of modules

- maybe even scan generated assembly for problematic instructions

Further work (replacing `GMP`):

- typed Assembly-level functions: `timesWord2# :: Word -> Word -> (# Word, Word #)`

- reduction of field elements scheduled via type-level lifted overflow information $\Rightarrow$ provably non-overflowing field elements at compile-time

- linear/affine type systems (Rust, Linear Haskell)
  - do not supplant CT-safety
  - may make more optimisations safe

- dependent types (Coq, Agda, Dependent Haskell)
  - state of the art approach used in implementation adopted by Mozilla Firefox
  - totality
  - type inference may be hard
  - efficient code generation is hard, so some subsets are used which mostly one-to-one generate C or Assembly constructs in a proof language

- linear/affine type systems (Rust, Linear Haskell)
  - do not supplant CT-safety
  - may make more optimisations safe
- dependent types (Coq, Agda, Dependent Haskell)
  - state of the art approach used in implementation adopted by Mozilla Firefox
  - totality
  - type inference may be hard
  - efficient code generation is hard, so some subsets are used which mostly one-to-one generate C or Assembly constructs in a proof language
    - proofs are not automatically transitive across compile chains due to different semantics, transitiveness is proven
    - but use of GCC (performance) vs. CompCert (verified)

Fin!