# Learning about Computer Arithmetic by Formally Verifying It

John Harrison

Amazon Web Services

RAIM Meeting 2025

Wed 5th Nov 2025 (11:00–12:00)

# 1998-2017: Verifying floating-point arithmetic at Intel



At the IEEE floating-point meeting 2006

# 2018-?: Verifying crypto bignums at AWS



AWS's Automated Reasoning Group in 2019

# Similarities

# Similarities

- ▶ When starting, largely ignorant about the subject matter

# Similarities

- When starting, largely ignorant about the subject matter
- Plenty of real problems motivate the application of verification

# Similarities

- ▶ When starting, largely ignorant about the subject matter
- ▶ Plenty of real problems motivate the application of verification
- ▶ Much information is/was 'folklore', not available in textbooks

# Similarities

- When starting, largely ignorant about the subject matter
- Plenty of real problems motivate the application of verification
- Much information is/was 'folklore', not available in textbooks
- At least we have precise specifications

# Similarities

- When starting, largely ignorant about the subject matter
- Plenty of real problems motivate the application of verification
- Much information is/was 'folklore', not available in textbooks
- At least we have precise specifications ... more or less

# Similarities

- When starting, largely ignorant about the subject matter
- Plenty of real problems motivate the application of verification
- Much information is/was 'folklore', not available in textbooks
- At least we have precise specifications . . . more or less
- Technological progress opens up new sub-areas of application

# Similarities

- When starting, largely ignorant about the subject matter
- Plenty of real problems motivate the application of verification
- Much information is/was 'folklore', not available in textbooks
- At least we have precise specifications . . . more or less
- Technological progress opens up new sub-areas of application
- Many mathematical similarities and analogies

# Differences

# Differences

- Open-sourcing
  - Intel work was internal and never made public
  - Cryptographic code an open-source library s2n-bignum

# Differences

- Open-sourcing
  - Intel work was internal and never made public
  - Cryptographic code an open-source library s2n-bignum
- Popularity/competition
  - Floating-point verification was relatively new when I started
  - Everyone and their dog did cryptographic FV even in 2018

# Differences

- ▶ Open-sourcing
  - ▶ Intel work was internal and never made public
  - ▶ Cryptographic code an open-source library s2n-bignum
- ▶ Popularity/competition
  - ▶ Floating-point verification was relatively new when I started
  - ▶ Everyone and their dog did cryptographic FV even in 2018
- ▶ Implementation and verification (much more satisfying)
  - ▶ Intel work was initially verifying existing code, but eventually led to new implementations
  - ▶ s2n-bignum was from the start implementation and verification together

# Differences

- Open-sourcing
  - Intel work was internal and never made public
  - Cryptographic code an open-source library s2n-bignum
- Popularity/competition
  - Floating-point verification was relatively new when I started
  - Everyone and their dog did cryptographic FV even in 2018
- Implementation and verification (much more satisfying)
  - Intel work was initially verifying existing code, but eventually led to new implementations
  - s2n-bignum was from the start implementation and verification together
- We want efficiency and correctness

# Differences

- ▶ Open-sourcing
  - ▶ Intel work was internal and never made public
  - ▶ Cryptographic code an open-source library s2n-bignum
- ▶ Popularity/competition
  - ▶ Floating-point verification was relatively new when I started
  - ▶ Everyone and their dog did cryptographic FV even in 2018
- ▶ Implementation and verification (much more satisfying)
  - ▶ Intel work was initially verifying existing code, but eventually led to new implementations
  - ▶ s2n-bignum was from the start implementation and verification together
- ▶ We want efficiency and correctness *and* security

# Differences

- Open-sourcing
    - Intel work was internal and never made public
    - Cryptographic code an open-source library s2n-bignum
- Popularity/competition
    - Floating-point verification was relatively new when I started
    - Everyone and their dog did cryptographic FV even in 2018
- Implementation and verification (much more satisfying)
    - Intel work was initially verifying existing code, but eventually led to new implementations
    - s2n-bignum was from the start implementation and verification together
- We want efficiency and correctness *and* security
    - 'Constant-time' code

# Motivations for FP verification



Intel's $475M mistake

# Motivations for crypto verification



https://nvd.nist.gov/vuln/detail/CVE-2017-3736

# ...and it's not just correctness



https://seclists.org/oss-sec/2018/q2/50

# THE PARIS256 ATTACK

*Or, Squeezing a Key Through a Carry Bit.*

Sean Devlin, Filippo Valsorda

## Introduction

We present an adaptive key recovery attack exploiting a small carry propagation bug in the Go standard library implementation of the NIST P-256 elliptic curve, reported to the Go project as issue 20040.

Following our attack, the vulnerability was assigned CVE-2017-8932, and caused the release of Go 1.7.6 and 1.8.2.

```
https://i.blackhat.com/us-18/Wed-August-8/
us-18-Valsorda-Squeezing-A-Key-Through-A-Carry-Bit-wp.
pdf
```

# From folklore to textbooks

# Things were very different when I started

Results were scattered in different papaers, with some tricks hardly written down at all:

# Things were very different when I started

Results were scattered in different papaers, with some tricks hardly written down at all:

- ▶ FMA-based division and square root: Markstein, *Computation of Elementary Functions on the IBM RISC System/6000 Processors*

# Things were very different when I started

Results were scattered in different papaers, with some tricks hardly written down at all:

- ▶ FMA-based division and square root: Markstein, *Computation of Elementary Functions on the IBM RISC System/6000 Processors*
- ▶ Computation of transcendental functions: Tang, *Table-driven implementation of the exponential function in IEEE floating-point arithmetic*

# Things were very different when I started

Results were scattered in different papaers, with some tricks hardly written down at all:

- ▶ FMA-based division and square root: Markstein, *Computation of Elementary Functions on the IBM RISC System/6000 Processors*
- ▶ Computation of transcendental functions: Tang, *Table-driven implementation of the exponential function in IEEE floating-point arithmetic*
- ▶ General floating-point "magic":
  - ▶ Sterbenz, *Floating-Point Computation*
  - ▶ Goldberg, *What every computer scientist should know about floating-point arithmetic*
  - ▶ Kahan, passim

# Exact sum and exact product

The exact sum property was relatively well-known (Sterbenz, Goldberg)

```
|- a IN iformat fmt ∧ b IN iformat fmt ∧
   a / 2 <= b ∧ b <= 2 * a
   ⇒ (b - a) IN iformat fmt
```

## Exact sum and exact product

The exact sum property was relatively well-known (Sterbenz, Goldberg)

```
|- a IN iformat fmt ∧ b IN iformat fmt ∧
   a / 2 <= b ∧ b <= 2 * a
   ⇒ (b - a) IN iformat fmt
```

The corresponding multiplicative one was more obscure, perhaps because FMA was then not widely used or standardized:

```
|- a IN iformat fmt ∧ b IN iformat fmt ∧
   2 pow (2 * precision fmt - 1) / 2 pow (ulpscale fmt)
   <= abs(a * b)
   ⇒ (a * b - round fmt Nearest (a * b)) IN iformat fmt
```

http://www.cs.berkeley.edu/~wkahan/ieee754status/
ieee754.ps

# We need the textbooks for cryptographic arithmetic!

These two together probably come closest.

# At least we have precise specifications . . .

# At least we have precise specifications . . .



For us, the spec is pretty easy to formalize, purely mathematical and almost formal already.

## . . . or do we?

It's customary to give a bound on the error in transcendental
functions in terms of 'units in the last place' (ulps), roughly the
gap between adjacent floating point numbers.

# ... or do we?

It's customary to give a bound on the error in transcendental functions in terms of 'units in the last place' (ulps), roughly the gap between adjacent floating point numbers.



At the boundary $2^k$ between 'binades', this distance changes, which makes it tricky.

# Goldberg's definition of ulp

*In general, if the floating-point number $d.d \cdots d \times \beta^e$ is used to represent $z$, it is in error by $|d.d \cdots d - (z/\beta^e)|\beta^{p-1}e$ units in the last place.*

# Goldberg's definition of ulp

*In general, if the floating-point number $d.d \cdots d \times \beta^e$ is used to represent $z$, it is in error by $|d.d \cdots d - (z/\beta^e)|\beta^{p-1}e$ units in the last place.*

So this is an error of $0.5 ulp$ according to Goldberg, but intuitively it should be $1 ulp$.

# Kahan's definition of ulp

$ulp(x)$ *is the gap between the two floating-point numbers nearest $x$, even if $x$ is one of them.*

# Kahan's definition of ulp

*ulp(x) is the gap between the two floating-point numbers nearest x, even if x is one of them.*

According to that definition this is an error of $0.4ulp$, but intuitively it should be $0.2ulp$. Rounding up is worse...

# Ambiguity in Ed25519 signature specification

This signature scheme is one of the most widely used for verifying the authenticity of messages, and is standardized in RFC 8032

```
          Edwards-Curve Digital Signature Algorithm (EdDSA)
```

```
Abstract

   This document describes elliptic curve signature scheme Edwards-curve
   Digital Signature Algorithm (EdDSA).  The algorithm is instantiated
   with recommended parameters for the edwards25519 and edwards448
   curves.  An example implementation and test vectors are provided.
```

# The definition of signature verification

The central operation in signature verification involves arithmetic on elliptic curve points $B$, $R$ and $A'$, with $[k]G$ denoting scalar multiplication of a group element $G$ by an integer $k$, i.e $G + G + \ldots + G$ ($k$ times)

# The definition of signature verification

The central operation in signature verification involves arithmetic
on elliptic curve points $B$, $R$ and $A'$, with $[k]G$ denoting scalar
multiplication of a group element $G$ by an integer $k$, i.e
$G + G + \ldots + G$ ($k$ times)

> 3. *Check the group equation* $[8][S]B = [8]R + [8][k]A'$.
> *It's sufficient, but not required, to instead check*
> $[S]B = R + [k]A'$.

# The definition of signature verification

The central operation in signature verification involves arithmetic on elliptic curve points $B$, $R$ and $A'$, with $[k]G$ denoting scalar multiplication of a group element $G$ by an integer $k$, i.e $G + G + \ldots + G$ ($k$ times)

> 3. Check the group equation $[8][S]B = [8]R + [8][k]A'$.
> It's sufficient, but not required, to instead check
> $[S]B = R + [k]A'$.

Most implementations check yet a different equation
$[S]B = R + [k \ \texttt{MOD} \ n]A$ where $n$ is the order of the basepoint $B$.

# The definition of signature verification

The central operation in signature verification involves arithmetic on elliptic curve points $B$, $R$ and $A'$, with $[k]G$ denoting scalar multiplication of a group element $G$ by an integer $k$, i.e $G + G + \ldots + G$ ($k$ times)

> 3. Check the group equation $[8][S]B = [8]R + [8][k]A'$. It's sufficient, but not required, to instead check $[S]B = R + [k]A'$.

Most implementations check yet a different equation $[S]B = R + [k \text{ MOD } n]A$ where $n$ is the order of the basepoint $B$. These are all equivalent for well-formed signatures and do not affect the key security properties. Nevertheless, the *full* group order is $8n$, so for general points on the curve, these three equations are inequivalent in general.

# New areas developing

As well as their 'classical' subject-matter, both floating-point arithmetic and cryptography are developing new emphases as a result of technological advances or market pressure:

# New areas developing

As well as their 'classical' subject-matter, both floating-point arithmetic and cryptography are developing new emphases as a result of technological advances or market pressure:

- ▶ Machine learning is leading to more emphasis on the standardization of low-precision floating-point numbers and renewing interest in alternative representations.

# New areas developing

As well as their 'classical' subject-matter, both floating-point arithmetic and cryptography are developing new emphases as a result of technological advances or market pressure:

▶ Machine learning is leading to more emphasis on the standardization of low-precision floating-point numbers and renewing interest in alternative representations.

▶ As well as classical RSA and elliptic curve techniques, post-quantum algorithms like the recently standardized ML-KEM and ML-DSA have lattice-based mathematical underpinnings.

# New areas developing

As well as their 'classical' subject-matter, both floating-point arithmetic and cryptography are developing new emphases as a result of technological advances or market pressure:

- ▶ Machine learning is leading to more emphasis on the standardization of low-precision floating-point numbers and renewing interest in alternative representations.
- ▶ As well as classical RSA and elliptic curve techniques, post-quantum algorithms like the recently standardized ML-KEM and ML-DSA have lattice-based mathematical underpinnings.

These new areas are rich and attractive targets for formal specification and verification.

# Mathematical similarity: bounds reasoning

In both domains, some formal counterpart of interval reasoning is useful, as implemented by tools like GAPPA:

# Mathematical similarity: bounds reasoning

In both domains, some formal counterpart of interval reasoning is useful, as implemented by tools like GAPPA:

- ▶ In floating-point arithmetic the end result may well be stated as a bound on some overall error term, and automation can help compute it.

# Mathematical similarity: bounds reasoning

In both domains, some formal counterpart of interval reasoning is useful, as implemented by tools like GAPPA:

- ▶ In floating-point arithmetic the end result may well be stated as a bound on some overall error term, and automation can help compute it.
- ▶ In *both* cases, it can be used to prove the absence of overflow/underflow so
    - ▶ Floating-point results stay finite and/or normalized giving better relative error biounds
    - ▶ Integer operations are exact because they don't wrap round.

# Mathematical analogy: MSB versus LSB algorithms

Floating-point numbers, usually being normalized, naturally lend themselves to algorithms based on the 'most significant bit'.

# Mathematical analogy: MSB versus LSB algorithms

Floating-point numbers, usually being normalized, naturally lend themselves to algorithms based on the 'most significant bit'.

In cryptography, even identifying and manipulating MSBs can be awkward or unnatural to do in constant time. More usage of LSB-based algorithms such as *Montgomery multiplication*.

# Mathematical analogy: MSB versus LSB algorithms

Floating-point numbers, usually being normalized, naturally lend themselves to algorithms based on the 'most significant bit'.

In cryptography, even identifying and manipulating MSBs can be awkward or unnatural to do in constant time. More usage of LSB-based algorithms such as *Montgomery multiplication*.

Nevertheless there are often strong similarities between 'metrical' (MSB) and '2-adic' (LSB) algorithms (see table in Brent-Zimmermann).

# Using Newton's method for reciprocals

Floating-point computation of $1/a$:

- ▶ Form initial approximation $y \approx \frac{1}{a}$
- ▶ Then iterate $y' = y \cdot (2 - ay) = y + y \cdot (1 - ay)$

# Using Newton's method for reciprocals

Floating-point computation of $1/a$:

- ▶ Form initial approximation $y \approx \frac{1}{a}$
- ▶ Then iterate $y' = y \cdot (2 - ay) = y + y \cdot (1 - ay)$

If $y = \frac{1}{a}(1 + \epsilon)$ then $y' = \frac{1}{a}(1 - \epsilon^2)$, the classic quadratic convergence where we get twice as many bits of accuracy per iteration.

# Modular inverses by Hensel lifting

Consider the 1-word (negated) modular inverse, called
`word_negmodinv` in s2n-bignum.

# Modular inverses by Hensel lifting

Consider the 1-word (negated) modular inverse, called
`word_negmodinv` in s2n-bignum.

Given a 64-bit unsigned and *odd* integer $a$, returns another 64-bit
integer $x$ such that $ax \equiv -1 \pmod{2^{64}}$, i.e. that
`a * x == 0xFFFFFFFFFFFFFFFF` using unsigned silently-wrapping
word operations like those on C's `uint64_t`.

# Modular inverses by Hensel lifting

Consider the 1-word (negated) modular inverse, called word_negmodinv in s2n-bignum.

Given a 64-bit unsigned and *odd* integer $a$, returns another 64-bit integer $x$ such that $ax \equiv -1 \pmod{2^{64}}$, i.e. that a * x == 0xFFFFFFFFFFFFFFFF using unsigned silently-wrapping word operations like those on C's uint64_t.

It is implemented in a directly similar way using Hensel lifting, the $p$-adic analog of Newton's method.

# Initial approximation

As with the floating-point inverse, we need an initial approximation to start with. The following piece of magic (in C syntax):

```
x = (a - (a<<2))^2
```

happens to give a 5-bit negated modular inverse, assuming a is odd.

# Hensel lifting step

Given a $k$-bit approximation $ax \equiv -1 \pmod{2^k}$, do the same Newton step with integers, except for a sign flip because we want a *negated* inverse:

```
e = a * x + 1;
y = e * x + x;
```

# Hensel lifting step

Given a $k$-bit approximation $ax \equiv -1 \pmod{2^k}$, do the same Newton step with integers, except for a sign flip because we want a *negated* inverse:

```
e = a * x + 1;
y = e * x + x;
```

By the initial assumption $ax = 2^k n - 1$ for some integer $n$

# Hensel lifting step

Given a $k$-bit approximation $ax \equiv -1 \pmod{2^k}$, do the same
Newton step with integers, except for a sign flip because we want a
*negated* inverse:

```
e = a * x + 1;
y = e * x + x;
```

By the initial assumption $ax = 2^k n - 1$ for some integer $n$
So $e = ax + 1 = 2^k n$ and

# Hensel lifting step

Given a $k$-bit approximation $ax \equiv -1 \pmod{2^k}$, do the same Newton step with integers, except for a sign flip because we want a *negated* inverse:

```
e = a * x + 1;
y = e * x + x;
```

By the initial assumption $ax = 2^k n - 1$ for some integer $n$
So $e = ax + 1 = 2^k n$ and $e + 1 = 2^k n + 1$

# Hensel lifting step

Given a $k$-bit approximation $ax \equiv -1 \pmod{2^k}$, do the same Newton step with integers, except for a sign flip because we want a *negated* inverse:

```
e = a * x + 1;
y = e * x + x;
```

By the initial assumption $ax = 2^k n - 1$ for some integer $n$
So $e = ax + 1 = 2^k n$ and $e + 1 = 2^k n + 1$ so then
$ay = ax(e + 1) = (2^k n - 1)(2^k n + 1) = 2^{2k} n^2 - 1$, i.e.
$ay \equiv -1 \pmod{2^{2k}}$.

# Open-source: s2n-bignum

A library of bignum arithmetic operations designed for
cryptographic applications.

# Open-source: s2n-bignum

A library of bignum arithmetic operations designed for cryptographic applications.

- ▶ Efficient: hand-crafted code with competitive performance

## Open-source: s2n-bignum

A library of bignum arithmetic operations designed for cryptographic applications.

▶ Efficient: hand-crafted code with competitive performance

▶ Correct: every function is formally verified mathematically

# Open-source: s2n-bignum

A library of bignum arithmetic operations designed for cryptographic applications.

- ▶ Efficient: hand-crafted code with competitive performance
- ▶ Correct: every function is formally verified mathematically
- ▶ Secure: all code is written in "constant-time" style

# Open-source: s2n-bignum

A library of bignum arithmetic operations designed for cryptographic applications.

- ▶ Efficient: hand-crafted code with competitive performance
- ▶ Correct: every function is formally verified mathematically
- ▶ Secure: all code is written in "constant-time" style

```
https://github.com/awslabs/s2n-bignum
```

# Open-source: s2n-bignum

A library of bignum arithmetic operations designed for cryptographic applications.

- ▶ Efficient: hand-crafted code with competitive performance
- ▶ Correct: every function is formally verified mathematically
- ▶ Secure: all code is written in "constant-time" style

```
https://github.com/awslabs/s2n-bignum
```

All hand-written or specially generated 64-bit ARM and x86 machine code.

# Coding and verification flow

# Comparison with ther crypto verification projects

Formally verified cryptography projects can be placed on this spectrum:

# Comparison with ther crypto verification projects

Formally verified cryptography projects can be placed on this spectrum:

- ▶ Correct-by-construction coding (HACL*, Jasmin)
- ▶ . . .
- ▶ A bit of both or somewhere in between (Fiat)
- ▶ . . .
- ▶ Separate verification (CryptoLine, s2n-bignum)

# Pros and cons of the s2n-bignum approach

# Pros and cons of the s2n-bignum approach

☺ Independent of compiler or even macro-assembler correctness.

# Pros and cons of the s2n-bignum approach

☺ Independent of compiler or even macro-assembler correctness.

☺ Applicable to highly tuned efficient code that is hard to generate automatically as well as compiled C code etc.

# Pros and cons of the s2n-bignum approach

☺ Independent of compiler or even macro-assembler correctness.

☺ Applicable to highly tuned efficient code that is hard to generate automatically as well as compiled C code etc.

☺ Code is conventional human-readable and human-modifiable, usage is independent of prover infrastructure.

# Pros and cons of the s2n-bignum approach

- ☺ Independent of compiler or even macro-assembler correctness.
- ☺ Applicable to highly tuned efficient code that is hard to generate automatically as well as compiled C code etc.
- ☺ Code is conventional human-readable and human-modifiable, usage is independent of prover infrastructure.
- ☹ Much more work involved writing code at this level, less structured representation.

# Pros and cons of the s2n-bignum approach

- ☺ Independent of compiler or even macro-assembler correctness.
- ☺ Applicable to highly tuned efficient code that is hard to generate automatically as well as compiled C code etc.
- ☺ Code is conventional human-readable and human-modifiable, usage is independent of prover infrastructure.
- ☹ Much more work involved writing code at this level, less structured representation.
- ☺/☹ Exposure of low-level details like exact stack and PC offsets and particular registers.

# Implementation and verification together

The Intel work started as verifying pre-existing code, but sometimes the proof unlocked efficiency improvements.

# Implementation and verification together

The Intel work started as verifying pre-existing code, but sometimes the proof unlocked efficiency improvements.

> *If $q$ is a floating point number within 1 ulp of the true quotient $a/b$ of two floating point numbers, and $y$ is the correctly rounded-to-nearest approximation of the exact reciprocal $\frac{1}{b}$, then the following iteration:*

$$
\begin{aligned}
r &= a - bq \\
q' &= q + ry
\end{aligned}
$$

> *using round-to-nearest in each case, yields the correctly rounded-to-nearest quotient $q'$. (Markstein)*

## Implementation and verification together

The Intel work started as verifying pre-existing code, but sometimes the proof unlocked efficiency improvements.

> If $q$ is a floating point number within 1 ulp of the true quotient $a/b$ of two floating point numbers, and $y$ is the correctly rounded-to-nearest approximation of the exact reciprocal $\frac{1}{b}$, then the following iteration:

$$
\begin{aligned}
r &= a - bq \\
q' &= q + ry
\end{aligned}
$$

> using round-to-nearest in each case, yields the correctly rounded-to-nearest quotient $q'$. (Markstein)

Simply changing to this unlocks more efficient algorithms:

> ... $y$ approximates $\frac{1}{b}$ to a relative error $< \frac{1}{2^p}$, ...

# Implementation and verification together

The Intel work started as verifying pre-existing code, but sometimes the proof unlocked efficiency improvements.

> *If $q$ is a floating point number within 1 ulp of the true quotient $a/b$ of two floating point numbers, and $y$ is the correctly rounded-to-nearest approximation of the exact reciprocal $\frac{1}{b}$, then the following iteration:*

$$
\begin{aligned}
r &= a - bq \\
q' &= q + ry
\end{aligned}
$$

> *using round-to-nearest in each case, yields the correctly rounded-to-nearest quotient $q'$. (Markstein)*

Simply changing to this unlocks more efficient algorithms:

> *... $y$ approximates $\frac{1}{b}$ to a relative error $< \frac{1}{2^p}$, ...*

*And* the proof is shorter!

# s2n-bignum code and proof

The s2n-bignum project was conceived from the beginning as integrated design and verification

# s2n-bignum code and proof

The s2n-bignum project was conceived from the beginning as integrated design and verification

- ▶ Each function has a precise specification and associated formal proof (no unverified code)

# s2n-bignum code and proof

The s2n-bignum project was conceived from the beginning as integrated design and verification

- ▶ Each function has a precise specification and associated formal proof (no unverified code)
- ▶ Code written with verification in mind, emphasizing simplicity where possible.

# s2n-bignum code and proof

The s2n-bignum project was conceived from the beginning as integrated design and verification

- ▶ Each function has a precise specification and associated formal proof (no unverified code)
- ▶ Code written with verification in mind, emphasizing simplicity where possible.
- ▶ Later we have sometimes taken code or algorithms from other libraries (e.g. Lenngrenn's X25519 code, `mlkem-native`) and added formal proofs post-hoc.

# s2n-bignum code and proof

The s2n-bignum project was conceived from the beginning as integrated design and verification

- ▶ Each function has a precise specification and associated formal proof (no unverified code)
- ▶ Code written with verification in mind, emphasizing simplicity where possible.
- ▶ Later we have sometimes taken code or algorithms from other libraries (e.g. Lenngrenn's X25519 code, `mlkem-native`) and added formal proofs post-hoc.
- ▶ Some of the algorithms would be difficult to trust *without* a formal proof.

# GCD and modular inverse using divstep

Classic binary gcd is attractive and simple

$$
\begin{aligned}
\gcd(2n, 2m) &= 2\gcd(n, m) \\
\gcd(2n+1, 2m) &= \gcd(2n+1, m) \\
\gcd(2n, 2m+1) &= gcd(n, 2m+1) \\
\gcd(2n+1, 2m+1) &= gcd(min(2n+1, 2m+1), |m-n|)
\end{aligned}
$$

# GCD and modular inverse using divstep

Classic binary gcd is attractive and simple

$$
\begin{aligned}
\gcd(2n, 2m) &= 2\gcd(n, m) \\
\gcd(2n + 1, 2m) &= \gcd(2n + 1, m) \\
\gcd(2n, 2m + 1) &= gcd(n, 2m + 1) \\
\gcd(2n + 1, 2m + 1) &= gcd(min(2n + 1, 2m + 1), |m - n|)
\end{aligned}
$$

The Bernstein-Yang divstep algorithm replaces the magnitude comparison with a single-word proxy $\delta$ (assume $n$ is odd):

$$
\mathsf{divstep}(\delta, n, m) = \begin{cases} (1 - \delta, m, (m - n)/2) & \text{if } \delta > 0 \wedge \mathsf{odd}(m) \\ (1 + \delta, n, (m + (m \bmod 2)n)/2) & \text{otherwise.} \end{cases}
$$

# GCD and modular inverse using divstep

Classic binary gcd is attractive and simple

$$
\begin{aligned}
\gcd(2n, 2m) &= 2\gcd(n, m) \\
\gcd(2n + 1, 2m) &= \gcd(2n + 1, m) \\
\gcd(2n, 2m + 1) &= gcd(n, 2m + 1) \\
\gcd(2n + 1, 2m + 1) &= gcd(min(2n + 1, 2m + 1), |m - n|)
\end{aligned}
$$

The Bernstein-Yang divstep algorithm replaces the magnitude comparison with a single-word proxy $\delta$ (assume $n$ is odd):

$$
\text{divstep}(\delta, n, m) = \begin{cases} (1 - \delta, m, (m - n)/2) & \text{if } \delta > 0 \wedge \text{odd}(m) \\ (1 + \delta, n, (m + (m \bmod 2)n)/2) & \text{otherwise.} \end{cases}
$$

However the bound/termination reasoning is *much* more complex.

# Dan Bernstein's formal bounds proof

Fortunately the inventor of the algorithm learned HOL Light and proved the bound for us:



← **Post**

**Daniel J. Bernstein** ✔
@hashbreaker

New formally verified proof of #safegcd iteration bound: cr.yp.to/2023/hull-ligh... (script for full run+extras: cr.yp.to/2023/hull-ligh...)
Advantages over previous formally verified proofs: (1) covers all input sizes; (2) finishes verifying in 10 minutes; (3) smaller TCB (HOL Light).

11:11 PM · Apr 16, 2023 · **14.3K** Views

💬 2          🔁 6          ♡ 30          🔖 3          ⬆

# Side-channels and "constant-time"

There are many side-channels by which systems may 'leak' secret info (like a private key) to an observer:

- ▶ Execution time
- ▶ Memory access pattern
- ▶ Power consumption
- ▶ Electromagnetic radiation emitted
- ▶ . . .
- ▶ Microarchitectural bugs

# Side-channels and "constant-time"

There are many side-channels by which systems may 'leak' secret info (like a private key) to an observer:

- Execution time ←
- Memory access pattern ←
- Power consumption
- Electromagnetic radiation emitted
- . . .
- Microarchitectural bugs

Main worries in typical multitasking OS on shared machine

# Side-channels and "constant-time"

There are many side-channels by which systems may 'leak' secret info (like a private key) to an observer:

- Execution time ←
- Memory access pattern ←
- Power consumption
- Electromagnetic radiation emitted
- . . .
- Microarchitectural bugs

Main worries in typical multitasking OS on shared machine

However, *Hertzbleed: Turning Power Side-Channel Attacks Into Remote Timing Attacks on x86* shows how to mount power attacks remotely via frequency scaling.

# How can you avoid timing/cache side-channels?

Want execution time, if not literally constant, *uncorrelated* with (secret) data being manipulated. How?

# How can you avoid timing/cache side-channels?

Want execution time, if not literally constant, *uncorrelated* with
(secret) data being manipulated. How?

- ▶ Add randomization or salting to the algorithm
- ▶ Balance timing of paths
- ▶ Just make it too fast to observe

# How can you avoid timing/cache side-channels?

Want execution time, if not literally constant, *uncorrelated* with (secret) data being manipulated. How?

- ▶ Add randomization or salting to the algorithm
- ▶ Balance timing of paths
- ▶ Just make it too fast to observe
- ▶ Always perform exactly the same operations regardless of (secret) data. ⟵ Our chosen solution

# How can you 'always do the same thing'?

When there is control flow depending on secret data:

```
if (n >= p) n = n - p;
```

# How can you 'always do the same thing'?

When there is control flow depending on secret data:

```
if (n >= p) n = n - p;
```

convert it into dataflow using masking, conditional moves etc.

```
b = (n < p) - 1;
n = n - (p & b);
```

## What about the compiler?

The compiler may naively turn mask creation back into a branch:

```
b = (n < p) - 1;
```

## What about the compiler?

The compiler may naively turn mask creation back into a branch:

```
b = (n < p) - 1;
```

Time to break out your copy of "Hacker's Delight":

```
b = (((~n & p) | ((~n | p) & (n - p))) >> 63) - 1;
```

# What about the compiler?

The compiler may naively turn mask creation back into a branch:

```
b = (n < p) - 1;
```

Time to break out your copy of "Hacker's Delight":

```
b = (((~n & p) | ((~n | p) & (n - p))) >> 63) - 1;
```

Another motivation for working directly in machine code where flags and useful instructions like CMOV and CSEL are available.

# Are the machine instructions constant-time?

- Some definitely not, e.g. division by zero is special
- General assumption that simple things like add, mul mostly are
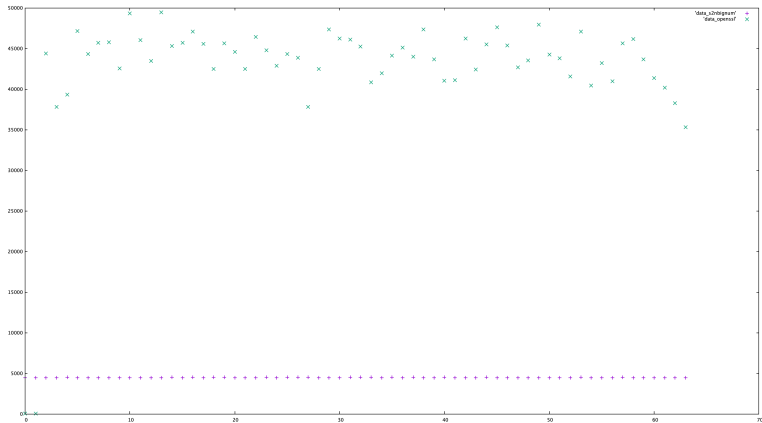
# Are the machine instructions constant-time?

- ▶ Some definitely not, e.g. division by zero is special
- ▶ General assumption that simple things like add, mul mostly are

Recently CPUs have started offering *some* guarantees (DIT bit or DOIT mode).

# Some empirical results on timing

Times for 384-bit modular inverse at bit densities 0–63,
nanoseconds on Intel® Xeon® Platinum 8175M, 2.5 GHz.

Questions?