# Performance Engineering of Proof-Based Software Systems at Scale

Jason Gross

Ph.D. Defense

MIT CSAIL

#### Takeaways

- Opportunity: Automate Verification to Enable Innovation
- Big Problem: Asymptotic Performance
- My Contribution: Reflective Partial Evaluation
- Important Next Steps



- Joint work with Andres Erbsen, Jade Philipoom, Adam Chlipala, et al
- Used in *majority* of secure connections from web browsers



ive.mozilla.org/foundation/identity-guidelines/firefox

HTTPS image modified from image by Sean MacEntee, CC BY 2.0, via Wikimedia Commons

November 30, 2020

Chrome logo from https://www.logo.wine/logo/Google Chrome ©2018 Google LLC All rights reserved. Chrome is a trademark of Google LLC. Go Logo By The Go Authors - https://blog.golang.org/go-brand, Public Domain, https://commons.wikimedia.org/w/index.php?curid=82371649 File:Logo of WireGuard.svg. (2020, April 21). Wikimedia Commons, the free media repository. Retrieved 23:20, November 28, 2020 from Wikimedia Commons Libra logo from By Libra Association - https://libra.org/, Public Domain, https://commons.wikimedia.org/w/index.php?curid=79808006

### Innovation with Cryptography

output[4] = r4;

"Better! Faster! Cheaper!"

- Hedging against more powerful attackers
- More mathematical security
- Reduce costs (server & user)



```
static inline void force_inline
fmul(felem output, const felem in2, const felem in) {
uint128_t t[5];
limb r0,r1,r2,r3,r4,s0,s1,s2,s3,s4,c;
r0 = in[0];
r1 = in[1];
r^2 = in[2]
r4 = in[4]
s\theta = in2[
s1 = in2[1]
<2 = in2[2];
s3 = in2[3];
s4 = in2[4];
t[0] = ((uint128_t) r0) * s0;
t[1] = ((uint128_t) r0) * s1 + ((uint128_t) r1) * s0;
t[2] = ((uint128_t) r0) * s2 + ((uint128_t) r2) * s0 + ((uint128_t) r1) * s1;
t[3] = ((uint128_t) r0) * s3 + ((uint128_t) r3) * s0 + ((uint128_t) r1) * s2 + ((uint128_t) r2) * s1;
t[4] = ((uint128 t) r0) * s4 + ((uint128 t) r4) * s0 + ((uint128 t) r3) * s1 + ((uint128 t) r1) * s3 + ((uint128 t) r2) * s2;
r4 *= 19;
r1 *= 19;
r2 *= 19;
r3 *= 19;
t[0] += ((uint128_t) r4) * s1 + ((uint128_t) r1) * s4 + ((uint128_t) r2) * s3 + ((uint128_t) r3) * s2;
t[1] += ((uint128_t) r4) * s2 + ((uint128_t) r2) * s4 + ((uint128_t) r3) * s3;
t[2] += ((uint128_t) r4) * s3 + ((uint128_t) r3) * s4;
t[3] += ((uint128_t) r4) * 54;
r0 = (limb)t[0] & 0x7ffffffffff; c = (limb)(t[0] >> 51);
t[1] += c; r1 = (limb)t[1] & 0x7fffffffffff; c = (limb)(t[1] >> 51);
t[2] += c; r2 = (limb)t[1] & 0x/fffffffff; c = (limb)(t[1] >> 51);
t[2] += c; r2 = (limb)t[2] & 0x/ffffffffff; c = (limb)(t[2] >> 51);
t[4] += c; r4 = (limb)t[4] & 0x/ffffffffff; c = (limb)(t[4] >> 51);
r0 += c * 19; c = r0 >> 51; r0 = r0 & 0x7ffffffffff;
r1 += c; c = r1 >> 51; r1 = r1 & 0x7fffffffffff;
r2 += c;
output[0] = r0;
output[1] = r1;
output[2] = r2;
output[3] = r3;
```

"Don't touch it; it works!"

- Lots of room for error
- Enormous cost of error
- Hard to find errors



#### The Overhead of Verification

• 10x—100x overhead

CompCert	5880	36,120	
seL4	8700	1,092,121	
CertiKOS	6500	96,642	
Fiat Cryptography	603	94,196	

Lines of Code	Lines of Verification
---------------	-----------------------

#### Automating Verification



November 30, 2020

Image modified from Thinking by ArmOkay from the Noun Project; script by Berkah Icon from the Noun Project; write by royyanandrian from the Noun Project

#### Our script is run and checked by...

- Dependently typed, interactive, tactic-driven proof assistants
- Dependently typed proof assistants are expressive
- Interactivity allows easy insertion of human ingenuity
- Tactics allow automation



#### The Big Problem in Automating Verification

#### • Asymptotic performance

- We can automate verification of toy examples in the proof engine
- BUT this automation takes way too long on real examples
- My work has been fixing this performance problem



#### The Potential of Automating Verification

Fiat Cryptography:



## It's really easy to use!



## It's really easy to use!



#### Requirements

1. Code we generate must be fast and constant time

Justification: server load, security

- 2. Easy to add and prove new algorithm, prime, architecture, ...
  - Justification: scalability of human effort, edit-compile-debug loops
- 3. Verification should not run forever

Justification: usability

#### Where was the asymptotic performance issue?

## Fiat Cryptography Pieces



Bounds Analysis



Bounds Analysis

#### Verification Time: 1 limb



#### Verification Time: 2 limbs



#### Verification Time: 3 limbs



Bounds Analysis

#### Partial Evaluation

November 30, 2020

#### What is Partial Evaluation?



Peanut Butter modified from image by P Thanga Vignesh from the Noun Project CC BY 3.0; Jam modified from image by Nikita Kozin from the Noun Project CC BY 3.0; Bread modified from image by Joanna Giansanti from the Noun Project CC BY 3.0; Sandwich bag by Kate Maldjian from the Noun Project; Shelves by Lluisa Iborra from the Noun Project; Refrigerator by shashank singh from the Noun Project

#### What is Partial Evaluation?

$$x + 2 + y - x + 6$$



#### Partial Evaluation in Fiat Cryptography

Evaluation

```
Template Code:
Definition mul (p q:list (Z*Z)):list (Z*Z) :=
                                                          Partial
Evaluation
 flat map (fun '(w, t) =>
   map (fun '(w', t') =>
      (w * w', t * t'))
    q) p.
Fixpoint square (p:list (Z*Z)):list (Z*Z)
:= match p with
    [] => []
    (w, t) :: ts
     => let two t := 2 * t in
         ((w * w, t * t))
           :: map (\lambda '(w', t'), (w * w', two_t * t')) ts)
         ++ square ts
  end.
Definition split (s:Z) (p:list(Z*Z)):list (Z*Z) * list (Z*Z)
:= let '(hi, lo) := partition (fun '(w, ) => w mod s =? 0) p in
```

```
(lo, map (fun '(w, t) => (w / s, t)) hi).
```

Definition reduce (s:Z) (c:list (Z\*Z)) (p:list (Z\*Z)):list (Z\*Z) := let '(lo, hi) := split s p in lo ++ mul c hi.

#### 64-bit square

Static Volu	11ac_23313_0100120 X2/	,
fiat_25519_carry_square(uint64	_t fiat_25519_uint128 x28	6
out1[5], const uint64_t arg1[5	<pre>]) fiat_25519_uint128 x29</pre>	;
{	fiat_25519_uint128 x30	15
<pre>uint64_t x1;</pre>	fiat_25519_uint128 x31	;
uint64_t x2;	<pre>uint64_t x32;</pre>	
uint64_t x3;	<pre>uint64_t x33;</pre>	
uint64_t x4;	fiat_25519_uint128 x34	5
uint64_t x5;	<pre>uint64_t x35;</pre>	
uint64_t x6;	<pre>uint64_t x36;</pre>	
uint64_t x7;	fiat_25519_uint128 x37	;
uint64_t x8;	<pre>uint64_t x38;</pre>	
fiat_25519_uint128 x9;	<pre>uint64_t x39;</pre>	
fiat_25519_uint128 x10;	fiat_25519_uint128 x40	;
fiat_25519_uint128 x11;	<pre>uint64_t x41;</pre>	
fiat_25519_uint128 x12;	<pre>uint64_t x42;</pre>	
fiat_25519_uint128 x13;	uint64_t x43;	
fiat_25519_uint128 x14;	uint64_t x44;	
fiat_25519_uint128 x15;	uint64_t x45;	
fiat_25519_uint128 x16;	uint64_t x46;	
fiat_25519_uint128 x17;	uint64_t x47;	
fiat_25519_uint128 x18;	fiat_25519_uint1 x48;	
fiat_25519_uint128 x19;	uint64_t x49;	
fiat_25519_uint128 x20;	uint64_t x50;	
fiat_25519_uint128 x21;	x1 = ((arg1[4]) *	
fiat_25519_uint128 x22;	UINT8_C(0x13));	
fiat_25519_uint128 x23;	x2 = (x1 * 0x2);	
fiat_25519_uint128 x24;	x3 = ((arg1[4]) * 0x2)	;
uint64_t x25;	x4 = ((arg1[3]) *	
<pre>uint64_t x26;</pre>	UINT8_C(0x13));	

#### 32-bit square

<pre>static void fiat_25519_carry_square(uint32_t</pre>	uint64_1
<pre>out1[10], const uint32_t arg1[10]) {</pre>	uint64_t
uint32_t x1;	uint64_t
uint32_t x2;	uint64_t
uint32_t x3;	uint64_t
uint32_t x4;	uint64_t
uint64_t x5;	uint64_t
uint32_t x6;	uint64_t
uint32_t x7;	uint64_t
uint32_t x8;	uint64_1
uint32_t x9;	uint64_t
uint32_t x10;	uint64_1
uint64_t x11;	uint64_t
uint32_t x12;	uint64_1
uint32_t x13;	uint64_t
uint32_t x14;	uint64_1
uint32_t x15;	uint64_1
uint32_t x16;	uint64_1
uint32_t x17;	uint64_1
uint32_t x18;	uint64_1
uint64_t x19;	uint64_t
uint64_t x20;	uint64_t
uint64_t x21;	uint64_1
uint64_t x22;	uint32_1
uint64_t x23;	uint64_1
uint64_t x24;	uint64_1
uint64_t x25;	uint64_1
uint64_t x26;	uint64_1
uint64_t x27;	uint64_1
uint64_t x28;	uint64_1
uint64_t x29;	uint64_1
uint64_t x30;	uint64_1
uint64_t x31;	uint64_1
uint64_t x32;	u1nt64_1
uint64_t x33;	u1nt64_1
uint64_t x34;	uint32_1
uint64_t x35;	u1nt64_1
uint64_t x36;	u1nt64_1
uinco4_c x57,	utilitisz_i
uint64_t x38;	u1nt64_1
uint64_t x39;	u1nt64_1
uinto4_t x40,	utilitiz_i
winted_t_x=1;	ulinto4_1
uinto4_c x42,	u1004_0
uinco4_c x+s,	utilitisz_i
wint64 t x45	uint64_1
winter a war.	uint 33
winted t v47:	uint64 4
wint64 t v49-	uint64_1
vieta e ven	uint 33
wint64 t v50-	uint64_1
winted t v51-	uint64_1
010004_C X31;	u111064_1

(5 = (x4 * 0x2);	x7);	<pre>UINT64_C(0x7ffffffffffff));</pre>
<pre>(6 = ((arg1[3]) * 0x2);</pre>	x18 =	x34 = (x32 + x29);
<pre>(7 = ((arg1[2]) * 0x2);</pre>	((fiat_25519_uint128)(arg1[1])	<pre>*x35 = (uint64_t)(x34 &gt;&gt; 51);</pre>
<pre>c8 = ((arg1[1]) * 0x2);</pre>	(arg1[1]));	x36 = (uint64_t)(x34 &
(9 =	x19 =	<pre>UINT64_C(0x7ffffffffffff));</pre>
(fiat_25519_uint128)(arg1[4])	*((fiat_25519_uint128)(arg1[0])	*x37 = (x35 + x28);
(1);	x3);	$x38 = (uint64_t)(x37 >> 51);$
(10 =	x20 =	$x39 = (uint64_t)(x37 \&$
(fiat_25519_uint128)(arg1[3])	*((fiat_25519_uint128)(arg1[0])	<pre>*UINT64_C(0x7ffffffffffff));</pre>
(2);	x6);	x40 = (x38 + x27);
(11 =	x21 =	$x41 = (uint64_t)(x40 >> 51);$
(fiat_25519_uint128)(arg1[3])	*((fiat_25519_uint128)(arg1[0])	*x42 = (uint64_t)(x40 &
(4);	x7);	<pre>UINT64_C(0x7ffffffffffff));</pre>
(12 =	x22 =	x43 = (x41 * UINT8_C(0x13));
(fiat_25519_uint128)(arg1[2])	*((fiat_25519_uint128)(arg1[0])	* x44 = (x26 + x43);
(2);	x8);	x45 = (x44 >> 51);
(13 =	x23 =	x46 = (x44 &
(fiat_25519_uint128)(arg1[2])	*((fiat_25519_uint128)(arg1[0])	*UINT64_C(0x7fffffffffff));
(5);	(arg1[0]));	x47 = (x45 + x33);
(14 =	x24 = (x23 + (x15 + x13));	x48 = (fiat_25519_uint1)(x47 >>
(fiat_25519_uint128)(arg1[2])	<pre>*x25 = (uint64_t)(x24 &gt;&gt; 51);</pre>	51);
[arg1[2]));	$x26 = (uint64_t)(x24 \&$	x49 = (x47 &
(15 =	UINT64_C(0x7fffffffffff));	<pre>UINT64_C(0x7fffffffffff));</pre>
(fiat_25519_uint128)(arg1[1])	*x27 = (x19 + (x16 + x14));	x50 = (x48 + x36);
(2);	x28 = (x20 + (x17 + x9));	out1[0] = x46;
(16 =	x29 = (x21 + (x18 + x10));	out1[1] = x49;
(fiat_25519_uint128)(arg1[1])	*x30 = (x22 + (x12 + x11));	out1[2] = x50;
(6);	x31 = (x25 + x30);	out1[3] = x39;
(17 =	<pre>x32 = (uint64_t)(x31 &gt;&gt; 51);</pre>	out1[4] = x42;
(fiat_25519_uint128)(arg1[1])	*x33 = (uint64_t)(x31 &	}

×107;	x41 = ((arg1[3]) * x5);	x88 = (uint32_t)(x86 & UINT32_C(0x1fffff));
x108;	x42 = ((arg1[3]) * ((uint64_t)x8 * 0x2));	x89 = (x87 + x84);
x109;	<pre>x43 = ((uint64_t)(arg1[3]) * x12);</pre>	x90 = (x89 >> 26);
×110;	<pre>x44 = ((uint64_t)(arg1[3]) * (x14 * 0x2));</pre>	<pre>x91 = (uint32_t)(x89 &amp; UINT32_C(0x3ffffff));</pre>
x111;	x45 = ((uint64_t)(arg1[3]) * x15);	x92 = (x90 + x83);
x112;	x46 = ((uint64_t)(arg1[3]) * ((arg1[3]) *	x93 = (x92 >> 25);
x113;	0x2));	<pre>x94 = (uint32_t)(x92 &amp; UINT32_C(0x1ffffff));</pre>
x114;	<pre>x47 = ((uint64_t)(arg1[2]) * x2);</pre>	x95 = (x93 + x82);
×115;	x48 = ((arg1[2]) * x5);	x96 = (x95 >> 26);
x116;	x49 = ((uint64_t)(arg1[2]) * x9);	<pre>x97 = (uint32_t)(x95 &amp; UINT32_C(0x3ffffff));</pre>
x117;	<pre>x50 = ((uint64_t)(arg1[2]) * x12);</pre>	x98 = (x96 + x81);
9_uint1 x118;	x51 = ((uint64_t)(arg1[2]) * x14);	x99 = (x98 >> 25);
×119;	x52 = ((uint64_t)(arg1[2]) * x15);	x100 = (uint32_t)(x98 &
x128;	x53 = ((uint64_t)(arg1[2]) * x16);	UINI32_C(0X1TTTTT));
<pre>gi[a]) * UINIS_c(ext3));</pre>	x54 = ((uint64_t)(argi[2]) * (argi[2]));	x181 = (x99 + x88);
- 0x2);	x55 = ((arg1[1]) * ((uint64_t)x2 * 0x2));	$x_{102} = (x_{101} >> 26);$
si[5]) * units ((0-12)).	$x_{30} = ((u_{10} + (u_{10} + (u_{$	$x_{105} = (010052 - 0)(x_{101} - a)$
gi[8]) * UINI8_C(0x13));	x57 = ((uint64_t)(arg1[1]) * (x9 * 0x2));	UINI32 ((0X3TTTTTT));
-1[0]) # 0-2);	$x_{56} = ((u_{10} + c_{10}) (arg1[1]) + x_{12}),$	$x_{104} = (x_{102} + x_{75}),$
si[3]) * unun ((2-12)).	$x_{33} = ((u_{11}c_{34} - c)(a_{1}g_{1}[1]) + (x_{14} - b_{22})),$	$x_{105} = (x_{104} + 3) (x_{104} + 8)$
* 0x2);	$x_{00} = ((u_{10} + 64 + )(a_{10} + 11) + (x_{10} + 6x_{2})),$	$x_{100} = (u_{10} x_{2} c)(x_{104} a)$
e1[7]) * 0x2):	x61 = ((uint64 ±)(ang1[1]) ± x17);	x107 = (x105 + x78):
ce1[6]) * UINTS C(0x13)):	$x_{63} = ((uint_{64} t)(arg1[1]) * ((arg1[1]) *$	$x_{100} = (x_{100} + x_{10}),$ $x_{108} = (x_{107} + x_{20} + 26);$
int64 t)x10 * 0x2):	9x2)):	$x109 = (uint32 t)(x107 \delta)$
cg1[6]) * 8x2):	$x_{54} = ((uint_{54} t)(arg1[0]) * x_3);$	UTNT32 C(0x3ffffff)):
re1[5]) * UINT8 C(0x13)):	$x_{65} = ((uint_{64} t)(arg1[0]) * x_{6});$	$x_{110} = (x_{100} + x_{77})$
rg1[5]) * 8x2):	$x_{66} = ((uint_{64} t)(arg1[0]) * x_{9});$	x111 = (x118 >> 25);
rg1[4]) * 0x2);	x67 = ((uint64 t)(arg1[0]) * x12);	x112 = (uint32 t)(x110 &
rg1[3]) * 0x2);	x68 = ((uint64 t)(arg1[0]) * x14);	UINT32 C(0x1ffffff)):
rg1[2]) * 0x2);	x69 = ((uint64 t)(arg1[0]) * x15);	x113 = (x111 * UINT8 C(0x13));
rg1[1]) * 0x2);	x70 = ((uint64 t)(arg1[0]) * x16);	x114 = (x76 + x113);
int64 t)(arg1[9]) * (x1 * 0x2));	<pre>x71 = ((uint64 t)(arg1[0]) * x17);</pre>	<pre>x115 = (uint32 t)(x114 &gt;&gt; 26);</pre>
int64_t)(arg1[8]) * x2);	<pre>x72 = ((uint64_t)(arg1[0]) * x18);</pre>	x116 = (uint32_t)(x114 &
<pre>int64_t)(arg1[8]) * x4);</pre>	<pre>x73 = ((uint64_t)(arg1[0]) * (arg1[0]));</pre>	UINT32_C(0x3ffffff));
rg1[7]) * ((uint64_t)x2 * 0x2));	x74 = (x73 + (x55 + (x48 + (x42 + (x37 +	x117 = (x115 + x88);
rg1[7]) * x5);	x33)))));	<pre>x118 = (fiat_25519_uint1)(x117 &gt;&gt; 25);</pre>
<pre>int64_t)(arg1[7]) * (x7 * 0x2));</pre>	x75 = (x74 >> 26);	<pre>x119 = (x117 &amp; UINT32_C(0x1ffffff));</pre>
<pre>int64_t)(arg1[6]) * x2);</pre>	<pre>x76 = (uint32_t)(x74 &amp; UINT32_C(0x3ffffff));</pre>	; x120 = (x118 + x91);
rg1[6]) * x5);	x77 = (x64 + (x56 + (x49 + (x43 + x38))));	out1[0] = x116;
<pre>int64_t)(arg1[6]) * x8);</pre>	x78 = (x65 + (x57 + (x50 + (x44 + (x39 +	out1[1] = x119;
<pre>int64_t)(arg1[6]) * x10);</pre>	x19)))));	out1[2] = x120;
rg1[5]) * ((uint64_t)x2 * 0x2));	x79 = (x66 + (x58 + (x51 + (x45 + x20))));	out1[3] = x94;
rg1[5]) * x5);	x80 = (x67 + (x59 + (x52 + (x46 + (x22 +	out1[4] = x97;
rg1[5]) * ((uint64_t)x8 * 0x2));	x21)))));	out1[5] = x100;
rg1[5]) * x11);	x81 = (x68 + (x68 + (x53 + (x25 + x23))));	out1[6] = x103;
int64_t)(arg1[5]) * (x13 * 0x2));	x82 = (x69 + (x61 + (x54 + (x29 + (x26 +	outi[/] = x106;
int64_t)(arg1[4]) * x2);	x24)))));	out1[8] = x109;
(SI(4)) · X3);	$x_{03} = (x_{10} + (x_{02} + (x_{34} + (x_{30} + x_{27})))))$	ouci[9] = x112,
		1
<pre>rg1(*j) · x11); int64 +)(nng1[4]) * x14).</pre>	$x_{20}$ (y7) + (y7) + (y71 + (y26 + y2)))).	
int64 t)(ang1[4]) t (ang1[4])):	vec = (v75 + v95).	
rg1[3]) * ((uint64 t)x2 * 8x2)):	x87 = (x85 >> 25);	

#### Partial Evaluation is slow



#### What is a proof engine?

- Declare a goal to prove
- Issue instructions to make partial progress on proving
- Can write scripts to automate issuing of instructions
- Tracks the progress and current state
- Can issue a trail (proof certificate) to be checked by a small checker ("kernel")

## Our Approach

- Dig deep to find the places of asymptotic blowup
- Understand the precise source of the blowup
- Fuse the different compiler passes deeply

#### Requirements for Partial Evaluation

•  $\beta$ -reduction

•  $\iota\delta$ -reduction + rewrites

code sharing preservation

#### $\beta$ -reduction

• Useful for eliminating function call overhead in the generated code, which is important for output code performance

• Example: ((λ x. x + 5) 2) ----> 2 + 5

#### $ι\delta$ -reduction + rewrites

- Useful for precomputation and eliminating function call overhead
- Arithmetic simplification necessary for getting right asymptotics of generated lines of code in fiat-crypto (quadratic vs. quartic)
- Example:

map (λ x. x + 5) [0; 1; z] ----> [(λ x. x + 5) 0; (λ x. x + 5) 1; (λ x. x + 5) z]

- Note that this leaves  $\boldsymbol{\beta}$  redexes
- Without  $\beta$ -reduction, this can blow up code size
- Fusing rewriting with  $\beta$ -reduction in a way that scales

#### Code Sharing Preservation

• Necessary for avoiding exponential blowup in generated code size

Example: map f (let y := x + x in let z := y + y in [z; z; z])
Iet y := x + x in let z := y + y in map f [z; z; z]
Iet y := x + x in let z := y + y in [f z; f z; f z]

• Fusing this with  $\beta$ - and  $\iota$ - reduction

## Compiler passes

- $\beta$ -reduction
  - eliminating function call overhead
- $\iota\delta$ -reduction + rewrites
  - inlining definitions to eliminate function call overhead
  - arithmetic simplification
- code sharing preservation
  - to avoid exponential blowup in code size

#### Extra Requirements

- Verified
  - Without extending the TCB
- Performant
  - Should not introduce extra super-linear factors

#### The compiler passes need to be fused

- Needed to achieve adequate asymptotic performance!
- Separating out rewriting results in quartic rather than quadratic loc
- Separating out ι-reduction (constant propagation) results in *enormous* code-size blowup
- Separating out code-sharing-preservation results in *enormous* codesize blowup

#### Implementation

- Reflective for performant and verified
- Normalization by Evaluation (NbE) (for β)
  - + let-lifting monad (code-sharing)
  - + rewriting ( $\iota\delta$ +rewrite)

#### Proof by Reflection

- Most steps in the proof engine make partial progress towards a goal and leave behind a trail
- Coq's proof engine has a highly optimized primitive step for validating the output of a computation
- Phrasing the goal so that we can just validate the output of a computation
  - Verifying the process, rather than having an ad-hoc process that leaves behind a trail verifying the output

#### Non-Reflection Example

Inductive is\_even:  $\mathbb{N} \rightarrow \mathbb{P}$  := |zero\_even : is\_even 0 |two\_plus\_even n : is\_even n  $\rightarrow$  is\_even (2+n).



November 30, 2020 Thinking modified from image by ArmOkay from the Noun Project CC BY 3.0; Coq logo from https://calebstanford.com/2019/01/15/coq-vector-image/ script by Berkah Icon from the Noun Project; write by royyanandrian from the Noun Project

#### Reflection Example: Up-Front Work

Inductive is\_even:  $\mathbb{N} \rightarrow \mathbb{P}$  := |zero\_even : is\_even 0 |two\_plus\_even n : is\_even n  $\rightarrow$  is\_even (2+n).

Inductive parity := even | odd.

Definition flip\_parity p
:= match p with even => odd | odd => even end.







Thinking modified from image by ArmOkay from the Noun Project CC BY 3.0; Coq logo from https://calebstanford.com/2019/01/15/coq-vector-image/

#### Reflection Example



Thinking modified from image by ArmOkay from the Noun Project CC BY 3.0 Coq logo from https://calebstanford.com/2019/01/15/coq-vector-image/

## Why reflective rewriting?

- Reflective rewriting is asymptotically faster
- The trail left by proof-engine-based rewriting is super-linear in the size of the code being transformed
- Tracking the goal incurs super-linear overhead in the number of binders
- Recursively computing only the output is asymptotically faster
- Side benefit: we can extract it to OCaml to run as a nifty commandline utility

#### Normalization by Evaluation

Goal: Reuse substitution in Gallina for substitution in ASTs Example: Turn " $(\lambda z p n x. z + (x + (p + n))) 0 1 (-1)$ " into  $(\lambda z p n x. (\lambda a b. rewrite("+", a, b))$  $z ((\lambda a b. rewrite("+", a, b)))$ x ( $(\lambda a b. rewrite("+", a, b))$ (rewrite("0")) (rewrite("1")) (rewrite("1")) Expression application -----> Gallina application ----> Gallina abstraction Expression abstraction  $\rightarrow$  rewriter invocations on  $\eta$ -expanded forms Expression constants

#### Normalization by Evaluation

Goal: Reuse substitution in Gallina for substitution in ASTs Example: Turn " $(\lambda z p n x. z + (x + (p + n))) 0 1 (-1)$ " into ( $\lambda z p n x. (\lambda a b. rewrite("+", a, b))$ z (( $\lambda a b. rewrite("+", a, b)$ ) x (( $\lambda a b. rewrite("+", a, b$ )) p n)))) (rewrite("0")) (rewrite("1")) (rewrite("-1")) Then reduce!

#### Normalization by Evaluation

Goal: Reuse substitution in Gallina for substitution in ASTs Example: Turn " $(\lambda z p n x. z + (x + (p + n))) 0 1 (-1)$ " into ( $\lambda z p n x. (\lambda a b. rewrite("+", a, b))$ z (( $\lambda a b. rewrite("+", a, b)$ ) x (( $\lambda a b. rewrite("+", a, b$ )) p n)))) (rewrite("0")) (rewrite("1")) (rewrite("-1")) Then reduce!

----> (λ x. rewrite("+", "0", rewrite("+", "x", rewrite("+", "1", "-1"))))

#### Let-Lifting

- Let-Lifting monad for code-sharing-preservation
- Assignment + return; bind is derived
- Rewrote NbE in this Let-Lifting monad
- Haven't seen it in the literature, but it's not too tricky
- Automatic ι-reduction was too tricky to figure out, so I hard-coded the cases we needed for fiat-crypto

#### Rewriting

- For  $\iota\delta$ +rewrite
- Using Parametric Higher-Order Abstract Syntax (PHOAS) to deal with binders allows delaying rewriting
- We thus achieve complete rewriting in a single pass when the rewrite rules form a DAG
  - We have extra magic for when they don't. The magic is called "fuel" and "try again".

#### More Features

- Select rewrite rule based on Coq's pattern matching so we don't need to walk the entire list of rewrite rules at every identifier/constant node just to see which ones apply
- On-the-fly emission of a type of codes for relevant constants
- Partial evaluation on the generated rewriter (further 2x efficiency)

#### Implementation

- Reflective for performant and verified
- Normalization by Evaluation (NbE) (for β)
  - + let-lifting monad (code-sharing)
  - + rewriting (ιδ+rewrite)

#### + more features

#### Evaluation

- It works!
- It's performant!



#### Performance



#### Performance on Fiat Cryptography



## Our Approach

- Dig deep to find the places of asymptotic blowup
- Understand the precise source of the blowup
- Fuse the different compiler passes deeply

#### Takeaways

- Opportunity: Automate Verification to Enable Innovation
- Big Problem: Asymptotic Performance
- My Contribution: Reflective Partial Evaluation
- Important Next Steps

#### Let's take a step back

- We succeeded, but this was very hard
- All of this to work around inadequate asymptotic performance of the proof engine
- This is typical!

#### What I did in my PhD



 Performance engineering (working around slowness in Coq)

Coding new things

Misc

#### Discovering interesting new things

### Our current approach to performance

- Using abstraction to prevent excessive unfolding
- Carving out the proof engine...
- ...and replacing it with reflection

#### Abstraction is not enough

- Systems code is often written in an adversarial context
- Symmetric crypto code is often written empirically
- Performant code breaks abstraction barriers

#### Reflection will not save us

- Using a proof assistant is for easily inserting human ingenuity to prove a broad range of things
- Using reflection is essentially giving up "easy" part
- As problems get bigger and harder and we need more ingenuity, it won't be cost-effective to do it reflectively
- Already in the partial evaluator I hit the same performance-scaling issues that I was trying to avoid by writing it in the first place (albeit at a smaller and surmountable scale)

## Can we avoid carving out the proof engine?

- Where is the performance issue?
- Turns out that it's pretty far from the problem we're solving
  - (This should be obvious, because if it wasn't, reflection wouldn't help.)
  - Example: evar instance allocation has nothing to do with correctness of a given C algorithm
- In my experience, it's not about generating a proof trail and it's not even really about individual steps being slow
  - It's about asymptotics of accessing and updating data being tracked
  - Sometimes just walking the term repeatedly is too much overhead

#### Not just an engineering challenge

- "Don't make stupid choices" isn't enough to get good asymptotic performance
- Try writing rewrite\_strat
  - inside the tactic engine
  - every step considered as progress towards proving something
  - linear in # of binders + # of rewrite locations + size of term
  - really hard, maybe impossible!
- We need to systematically study proof engines with an eye towards asymptotic performance!

#### Next Questions about Proof Engines

- Where does the performance overhead really come from?
- What things are people not currently doing due to performance overhead?
- What is an adequate set of primitives?
- What are acceptable thresholds on asymptotic behavior?
- Is it possible to achieve adequate performance simultaneously on all the primitives?



I think solving this problem—getting the basics of proof engines right, asymptotically—will drastically accelerate the scale of what we as a field can handle, and bring verification closer to its promise and potential of enabling innovation in industry.

# Thank you for your time and attention!

# Questions?