

# High-Throughput Elliptic Curve Cryptography using AVX2 Vector Instructions

---

Hao Cheng   Johann Großschädl   Jiaqi Tian   Peter B. Rønne   Peter Y. A. Ryan

University of Luxembourg

SAC 2020



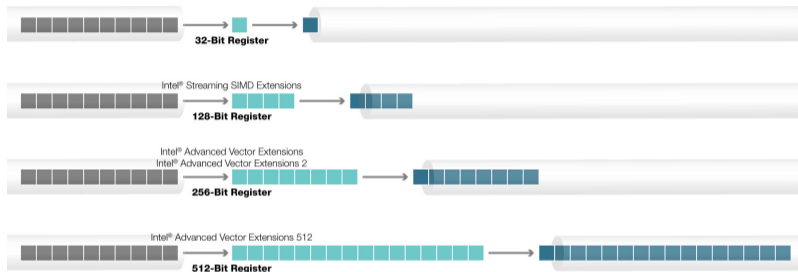
# SIMD



Single Instruction **M**ultiple **D**ata

# Intel x86/x64 Vector Extensions

ISA	Year	SIMD registers		
		#	FP	Integer
MMX	1997	8	64-bit	64-bit
SSE	1999	8	128-bit	128-bit
AVX	2011	16	256-bit	128-bit
AVX2	2013	16	256-bit	256-bit
AVX512	2016	32	512-bit	512-bit

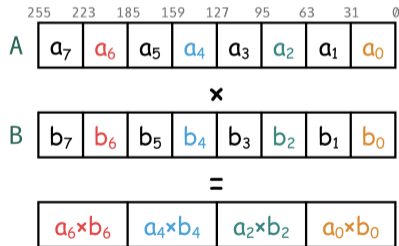


Intel® Advanced Vector eXtensions (AVX) series<sup>1</sup> (bottom two rows)

<sup>1</sup>figure from <https://www.prowesscorp.com/what-is-intel-avx-512-and-why-does-it-matter/>

## Properties

- ⊙ SIMD fashions : 8-bit × 32    16-bit × 16    32-bit × 8    64-bit × 4
- ⊙ Multiplier : 32-bit



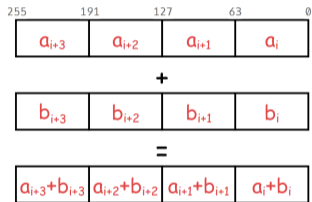
```
__m256i _mm256_mul_epu32 (__m256i A, __m256i B)
```

## ECC with SIMD Acceleration

- 1) Field arithmetic ← limbs
- 2) Curve arithmetic ← field operations
- 3) Combination of ← 1) and 2)
- 4) Mixed use of 1), 2) and 3)

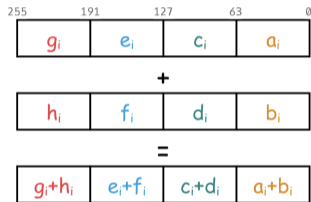
# ECC with SIMD Acceleration

- 1) Field arithmetic ← limbs
- 2) Curve arithmetic ← field operations
- 3) Combination of 1) and 2) ← 1) and 2)
- 4) Mixed use of 1), 2) and 3)



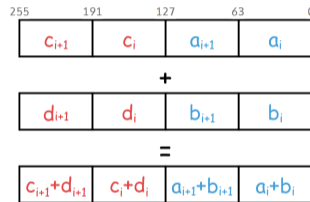
1) Accelerate the single addition

$$A + B$$



2) Parallel additions

$$G + H \mid E + F \mid C + D \mid A + B$$



3) Combination of 1) and 2)

$$C + D \mid A + B$$

( Each  $x_i$  is one limb of the large integer  $X$  )

# $(n \times m)$ -Way Parallelism

the number of  
field operations



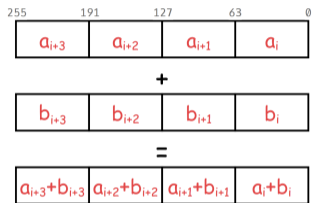
$(n \times m)$ -way



the number of elements  
used by each field operation

# $(n \times m)$ -Way Parallelism

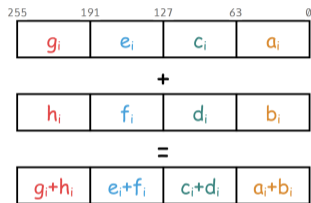
$(1 \times 4)$ -way



1) Accelerate the single addition

$$A + B$$

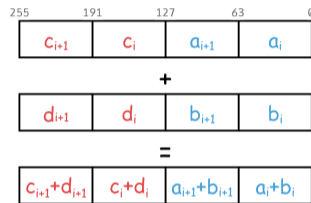
$(4 \times 1)$ -way



2) Parallel additions

$$G + H \mid E + F \mid C + D \mid A + B$$

$(2 \times 2)$ -way



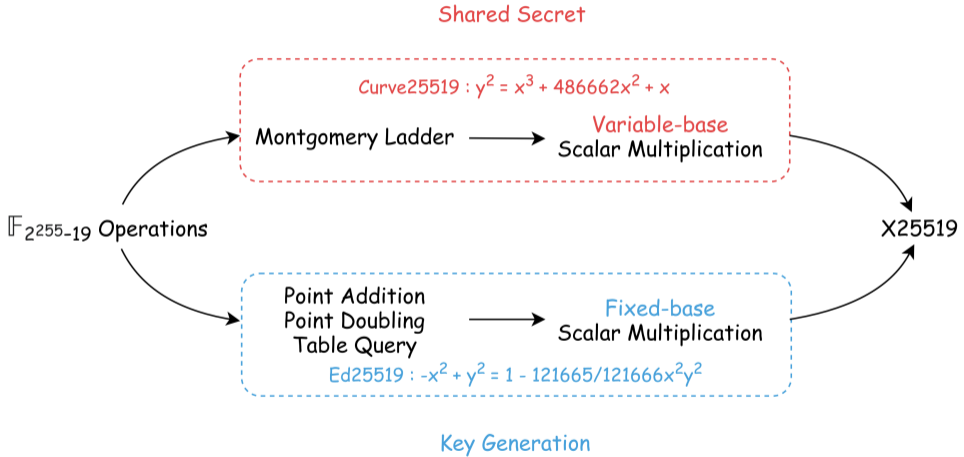
3) Combination of 1) and 2)

$$C + D \mid A + B$$

( Each  $x_i$  is one limb of the large integer X )



# X25519



# Latency-Optimized Work

## Low-Latency X25519 using AVX

Work	Authors	ISA	Impl.	Var-base scalar mul.
[Chou15]	Chou	AVX	(2 × 1)-way	137.2 k cycles
[FHLD19]	Faz-H., López, Dahab	AVX2	(2 × 2)-way	99.4 k cycles
[HEY20]	Hisil, Egrice, Yassi	AVX512	(4 × 2)-way	74.4 k cycles
[NS20]	Nath, Sarkar	AVX2 assembly	(4 × 1)-way	95.4 k cycles

# Latency-Optimized Work

## Low-Latency X25519 using AVX

Work	Authors	ISA	Impl.	Var-base scalar mul.
[Chou15]	Chou	AVX	(2 × 1)-way	137.2 k cycles <b>27.6%</b>
[FHLD19]	Faz-H., López, Dahab	AVX2	(2 × 2)-way	99.4 k cycles
[HEY20]	Hisil, Egrice, Yassi	AVX512	(4 × 2)-way	74.4 k cycles
[NS20]	Nath, Sarkar	AVX2 assembly	(4 × 1)-way	95.4 k cycles

# Latency-Optimized Work

## Low-Latency X25519 using AVX

Work	Authors	ISA	Impl.	Var-base scalar mul.	
[Chou15]	Chou	AVX	(2 × 1)-way	137.2 k cycles	27.6%
[FHLD19]	Faz-H., López, Dahab	AVX2	(2 × 2)-way	99.4 k cycles	25.2%
[HEY20]	Hisil, Egrice, Yassi	AVX512	(4 × 2)-way	74.4 k cycles	
[NS20]	Nath, Sarkar	AVX2 assembly	(4 × 1)-way	95.4 k cycles	

# Latency-Optimized Work

## Low-Latency X25519 using AVX

Work	Authors	ISA	Impl.	Var-base scalar mul.	
[Chou15]	Chou	AVX	(2 × 1)-way	137.2 k cycles	27.6%
[FHLD19]	Faz-H., López, Dahab	AVX2	(2 × 2)-way	99.4 k cycles	25.2%
[HEY20]	Hisil, Egrice, Yassi	AVX512	(4 × 2)-way	74.4 k cycles	
[NS20]	Nath, Sarkar	AVX2 assembly	(4 × 1)-way	95.4 k cycles	

Do not scale very well!



# Throughput v.s. Latency

- ⦿ How to exploit the massive parallelism of future SIMD extensions?

# Throughput v.s. Latency

- ⊙ How to exploit the massive parallelism of future SIMD extensions?
- ⊙ Why **low-latency** implementations?

# Throughput v.s. Latency

- ⊙ How to exploit the massive parallelism of future SIMD extensions?
- ⊙ Why **low-latency** implementations?
  - reduces the overall handshake-latency for a TLS client side



# Throughput v.s. Latency

- ⊙ How to exploit the massive parallelism of future SIMD extensions?
- ⊙ Why **low-latency** implementations?
  - reduces the overall handshake-latency for a TLS client side

Computation  $\ll$  Transmission !

# Throughput v.s. Latency

- ⊙ How to exploit the massive parallelism of future SIMD extensions?
- ⊙ Why **low-latency** implementations?
  - reduces the overall handshake-latency for a TLS client side

Computation  $\ll$  Transmission !

- ⊙ Why **high-throughput** implementations?

# Throughput v.s. Latency

## Why throughput-optimized?

TLS servers of big organizations ← several 10,000 TLS handshakes per second

- Latency ✗
- Throughput ✓

# Throughput v.s. Latency

## Why throughput-optimized?

TLS servers of big organizations ← several 10,000 TLS handshakes per second

- Latency ✗
- Throughput ✓

High throughput instead of low latency?

What throughput can it achieve?

# This Work

- ⦿ Takes first step to answer these questions
- ⦿ Introduces a throughput-optimized AVX2 implementation of X25519
  - **variable-base** scalar multiplication on **Curve25519**
  - **fixed-base** scalar multiplication on **Ed25519**



## Methodology – (4 × 1)-way scalar multiplication

Perform **FOUR** scalar multiplications simultaneously!

## “Coarse-Grained” Parallelism

- ⊙ Scalar multiplication
  - ⊙ Point arithmetic
  - ⊙ Field arithmetic
- } 64-bit element of 256-bit AVX2 vector

# “Coarse-Grained” Parallelism

- ⊙ Scalar multiplication
  - ⊙ Point arithmetic
  - ⊙ Field arithmetic
- } 64-bit element of 256-bit AVX2 vector

## Advantages

- 1) Easy to implement
- 2) Fully exploit parallelism
- 3) Support various SIMD extensions (straightforward extension to AVX512)



# Multi-Precision Representation

Radix- $2^{25.5}$  (e.g. [Chou15], [FHLD19])

$$f = f_0 + 2^{26} f_1 + 2^{51} f_2 + 2^{77} f_3 + 2^{102} f_4 + 2^{128} f_5 + 2^{153} f_6 + 2^{179} f_7 + 2^{204} f_8 + 2^{230} f_9$$

# Multi-Precision Representation

Radix- $2^{25.5}$  (e.g. [Chou15], [FHLD19])

$$f = f_0 + 2^{26} f_1 + 2^{51} f_2 + 2^{77} f_3 + 2^{102} f_4 + 2^{128} f_5 + 2^{153} f_6 + 2^{179} f_7 + 2^{204} f_8 + 2^{230} f_9$$

Radix- $2^{29}$  (this work)

$$f = f_0 + 2^{29} f_1 + 2^{58} f_2 + 2^{87} f_3 + 2^{116} f_4 + 2^{145} f_5 + 2^{174} f_6 + 2^{203} f_7 + 2^{232} f_8$$

# Multi-Precision Representation

Radix- $2^{25.5}$  (e.g. [Chou15], [FHLD19])

$$f = f_0 + 2^{26} f_1 + 2^{51} f_2 + 2^{77} f_3 + 2^{102} f_4 + 2^{128} f_5 + 2^{153} f_6 + 2^{179} f_7 + 2^{204} f_8 + 2^{230} f_9$$

Radix- $2^{29}$  (this work)

$$f = f_0 + 2^{29} f_1 + 2^{58} f_2 + 2^{87} f_3 + 2^{116} f_4 + 2^{145} f_5 + 2^{174} f_6 + 2^{203} f_7 + 2^{232} f_8$$

⊙  $(2 \times 2)$ -way    both use five limbs

# Multi-Precision Representation

Radix- $2^{25.5}$  (e.g. [Chou15], [FHLD19])

$$f = f_0 + 2^{26} f_1 + 2^{51} f_2 + 2^{77} f_3 + 2^{102} f_4 + 2^{128} f_5 + 2^{153} f_6 + 2^{179} f_7 + 2^{204} f_8 + 2^{230} f_9$$

Radix- $2^{29}$  (this work)

$$f = f_0 + 2^{29} f_1 + 2^{58} f_2 + 2^{87} f_3 + 2^{116} f_4 + 2^{145} f_5 + 2^{174} f_6 + 2^{203} f_7 + 2^{232} f_8$$

- ⊙  $(2 \times 2)$ -way      both use five limbs
- ⊙  $(4 \times 1)$ -way      Radix- $2^{25.5}$  uses ten limbs  
                         Radix- $2^{29}$  uses nine limbs

# Field Element Vector Set

$$\begin{array}{cccccccccccc}
 h = & 2^0 & \color{red}{h_0} & +2^{29} & \color{red}{h_1} & +2^{58} & \color{red}{h_2} & +2^{87} & \color{red}{h_3} & +2^{116} & \color{red}{h_4} & +2^{145} & \color{red}{h_5} & +2^{174} & \color{red}{h_6} & +2^{203} & \color{red}{h_7} & +2^{232} & \color{red}{h_8} & \begin{array}{l} 220 \\ 192 \end{array} \\
 g = & 2^0 & \color{blue}{g_0} & +2^{29} & \color{blue}{g_1} & +2^{58} & \color{blue}{g_2} & +2^{87} & \color{blue}{g_3} & +2^{116} & \color{blue}{g_4} & +2^{145} & \color{blue}{g_5} & +2^{174} & \color{blue}{g_6} & +2^{203} & \color{blue}{g_7} & +2^{232} & \color{blue}{g_8} & \begin{array}{l} 156 \\ 128 \end{array} \\
 f = & 2^0 & \color{teal}{f_0} & +2^{29} & \color{teal}{f_1} & +2^{58} & \color{teal}{f_2} & +2^{87} & \color{teal}{f_3} & +2^{116} & \color{teal}{f_4} & +2^{145} & \color{teal}{f_5} & +2^{174} & \color{teal}{f_6} & +2^{203} & \color{teal}{f_7} & +2^{232} & \color{teal}{f_8} & \begin{array}{l} 92 \\ 64 \end{array} \\
 e = & 2^0 & \color{orange}{e_0} & +2^{29} & \color{orange}{e_1} & +2^{58} & \color{orange}{e_2} & +2^{87} & \color{orange}{e_3} & +2^{116} & \color{orange}{e_4} & +2^{145} & \color{orange}{e_5} & +2^{174} & \color{orange}{e_6} & +2^{203} & \color{orange}{e_7} & +2^{232} & \color{orange}{e_8} & \begin{array}{l} 28 \\ 0 \end{array} \\
 \hline
 \mathbf{A} = & 2^0 & \mathbf{a}_0 & +2^{29} & \mathbf{a}_1 & +2^{58} & \mathbf{a}_2 & +2^{87} & \mathbf{a}_3 & +2^{116} & \mathbf{a}_4 & +2^{145} & \mathbf{a}_5 & +2^{174} & \mathbf{a}_6 & +2^{203} & \mathbf{a}_7 & +2^{232} & \mathbf{a}_8 & 
 \end{array}$$

$$\begin{aligned}
 \mathbf{A} = [e, f, g, h] &= \left[ \sum_{i=0}^8 2^{29i} e_i, \sum_{i=0}^8 2^{29i} f_i, \sum_{i=0}^8 2^{29i} g_i, \sum_{i=0}^8 2^{29i} h_i \right] \\
 &= \sum_{i=0}^8 2^{29i} [e_i, f_i, g_i, h_i] = \sum_{i=0}^8 2^{29i} \mathbf{a}_i \quad \text{with} \quad \mathbf{a}_i = [e_i, f_i, g_i, h_i].
 \end{aligned}$$

## Arithmetic in $\mathbb{F}_{2^{255}-19}$

⊙ Modulus  $p = 2^6 \cdot (2^{255} - 19) \leftarrow 29\text{-bit} \times 9 = 261\text{-bit}$

## Arithmetic in $\mathbb{F}_{2^{255}-19}$

- ⊙ Modulus  $p = 2^6 \cdot (2^{255} - 19) \leftarrow 29\text{-bit} \times 9 = 261\text{-bit}$
- ⊙ Addition  $\rightarrow$  ordinary integer addition  $r = a + b$

## Arithmetic in $\mathbb{F}_{2^{255}-19}$

- ⊙ Modulus  $p = 2^6 \cdot (2^{255} - 19) \leftarrow 29\text{-bit} \times 9 = 261\text{-bit}$
- ⊙ Addition  $\rightarrow$  ordinary integer addition  $r = a + b$
- ⊙ Subtraction
  - ordinary subtraction  $r = 2p + a - b$
  - modular subtraction  $r = 2p + a - b \bmod p$



## Arithmetic in $\mathbb{F}_{2^{255}-19}$

- ⊙ Modulus  $p = 2^6 \cdot (2^{255} - 19) \leftarrow 29\text{-bit} \times 9 = 261\text{-bit}$
- ⊙ Addition  $\rightarrow$  ordinary integer addition  $r = a + b$
- ⊙ Subtraction
  - ordinary subtraction  $r = 2p + a - b$
  - modular subtraction  $r = 2p + a - b \bmod p$
- ⊙ **Multiplication**  $r = a * b \bmod p$

## Arithmetic in $\mathbb{F}_{2^{255}-19}$

- ⊙ Modulus  $p = 2^6 \cdot (2^{255} - 19) \leftarrow 29\text{-bit} \times 9 = 261\text{-bit}$
- ⊙ Addition  $\rightarrow$  ordinary integer addition  $r = a + b$
- ⊙ Subtraction
  - ordinary subtraction  $r = 2p + a - b$
  - modular subtraction  $r = 2p + a - b \bmod p$
- ⊙ **Multiplication**  $r = a * b \bmod p$
- ⊙ Squaring  $\rightarrow$  special multiplication  $r = a^2 = a * a \bmod p$

# Field Multiplication

## Design Principles

- ⦿ Make full use of execution ports
- ⦿ Reduce the sequential dependencies

# Field Multiplication

## Design Principles

- ⦿ Make full use of execution ports
- ⦿ Reduce the sequential dependencies

a dozen of candidates → benchmark → the lowest latency



# Field Multiplication

## Design Principles

- ⦿ Make full use of execution ports
- ⦿ Reduce the sequential dependencies

a dozen of candidates → benchmark → the lowest latency

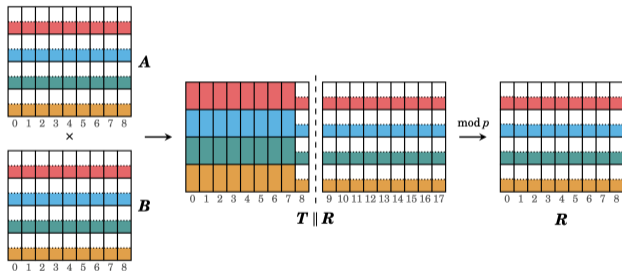


## Distinctions of candidates

- 1) Reduction & multiplication → **separated** or **interleaved**?
- 2) Different carry propagation plans
- 3) Intermediate values → local variables?

# Field Multiplication

```
1 #include <immintrin.h>
2 #define ADD(X,Y) _mm256_add_epi64(X,Y) /* VPADDQ */
3 #define MUL(X,Y) _mm256_mul_epu32(X,Y) /* VPMULUDQ */
4 #define AND(X,Y) _mm256_and_si256(X,Y) /* VPAND */
5 #define SRL(X,Y) _mm256_srli_epi64(X,Y) /* VPSRLQ */
6 #define BCAST(X) _mm256_set1_epi64x(X) /* VPBROADCASTQ */
7 #define MASK29 0xffffffff /* mask of 29 LSBs */
8
9 void fp_mul(__m256i *r, const __m256i *a, const __m256i *b)
10 {
11     int i, j, k; __m256i t[9], accu;
12
13     /* 1st loop of the product-scanning multiplication */
14     for (i = 0; i < 9; i++) {
15         t[i] = BCAST(0);
16         for (j = 0, k = i; k >= 0; j++, k--)
17             t[i] = ADD(t[i], MUL(a[j], b[k]));
18     }
19     accu = SRL(t[8], 29);
20     t[8] = AND(t[8], BCAST(MASK29));
21
22     /* 2nd loop of the product-scanning multiplication */
23     for (i = 9; i < 17; i++) {
24         for (j = i-8, k = 8; j < 9; j++, k--)
25             accu = ADD(accu, MUL(a[j], b[k]));
26         r[i-9] = AND(accu, BCAST(MASK29));
27         accu = SRL(accu, 29);
28     }
29     r[8] = accu;
30
31     /* modulo reduction and conversion to 29-bit limbs */
32     accu = BCAST(0);
33     for (i = 0; i < 9; i++) {
34         accu = ADD(accu, MUL(r[i], BCAST(64*19)));
35         accu = ADD(accu, t[i]);
36         r[i] = AND(accu, BCAST(MASK29));
37         accu = SRL(accu, 29);
38     }
39
40     /* limbs in r[0] can finally be 30 bits long */
41     r[0] = ADD(r[0], MUL(accu, BCAST(64*19)));
42 }
```



# Point Arithmetic

Take advantage of two types of field subtraction !

```
1 void point_add(ExtPoint *R, ExtPoint *P, ProPoint *Q)
2 {
3     __m256i t[9];
4
5     fp_mul(t, P->e, P->h);           /*  $T = E_{\mathcal{P}} \times H_{\mathcal{P}}$  */
6     fp_sub(R->e, P->y, P->x);         /*  $E_{\mathcal{R}} = Y_{\mathcal{P}} - X_{\mathcal{P}}$  */
7     fp_add(R->h, P->y, P->x);         /*  $H_{\mathcal{R}} = Y_{\mathcal{P}} + X_{\mathcal{P}}$  */
8     fp_mul(R->x, R->e, Q->y);         /*  $X_{\mathcal{R}} = E_{\mathcal{R}} \times Y_{\mathcal{Q}}$  */
9     fp_mul(R->y, R->h, Q->x);         /*  $Y_{\mathcal{R}} = H_{\mathcal{R}} \times X_{\mathcal{Q}}$  */
10    fp_sub(R->e, R->y, R->x);          /*  $E_{\mathcal{R}} = Y_{\mathcal{R}} - X_{\mathcal{R}}$  */
11    fp_add(R->h, R->y, R->x);          /*  $H_{\mathcal{R}} = Y_{\mathcal{R}} + X_{\mathcal{R}}$  */
12    fp_mul(R->x, t, Q->z);             /*  $X_{\mathcal{R}} = T \times Z_{\mathcal{Q}}$  */
13    fp_sbc(t, P->z, R->x);             /*  $T = Z_{\mathcal{P}} - X_{\mathcal{R}}$  */
14    fp_add(R->x, P->z, R->x);          /*  $X_{\mathcal{R}} = Z_{\mathcal{P}} + X_{\mathcal{R}}$  */
15    fp_mul(R->z, t, R->x);             /*  $Z_{\mathcal{R}} = T \times X_{\mathcal{R}}$  */
16    fp_mul(R->y, R->x, R->h);          /*  $Y_{\mathcal{R}} = X_{\mathcal{R}} \times H_{\mathcal{R}}$  */
17    fp_mul(R->x, R->e, t);             /*  $X_{\mathcal{R}} = E_{\mathcal{R}} \times T$  */
18 }
```

# Measurement Environment

## Platform

- a **Haswell** Intel® Core™ i7-4710HQ CPU clocked at 2.5 GHz
- a **Skylake** Intel® Core™ i5-6360U CPU clocked at 2.0 GHz

⊙ Compiler Clang 10.0.0

⊙ Disabled Features

- Intel® Turbo Boost ✗
- Intel® Hyper-Threading ✗



# Performance Evaluation

CPU cycles of  $(4 \times 1)$ -way field and point arithmetic

Domain	Operation	[FHLD19]		This Work	
		Haswell	Skylake	Haswell	Skylake
$\mathbb{F}_{2^{255}-19}$	Addition	12	12	11	11
	Ord. Subtraction	n/a	n/a	14	12
	Mod. Subtraction	n/a	n/a	32	31
	Multiplication	159	105	<b>122</b>	<b>88</b>
	Squaring	114	85	<b>87</b>	<b>65</b>
twisted Edwards curve	Point Addition	1096	833	965	705
	Point Doubling	n/a	n/a	830	624
	Table Query	208	201	218	205
Montgomery curve	Ladder Step	n/a	n/a	1118	818

# Performance Evaluation

Platform	CPU Frequency	Key Generation		Shared Secret		Table Size
		Latency	Throughput	Latency	Throughput	
Haswell	2.5 GHz	104,579 cycles	95,568 ops/sec	329,455 cycles	30,336 ops/sec	24 kB
Skylake	2.0 GHz	80,249 cycles	99,363 ops/sec	246,636 cycles	32,318 ops/sec	24 kB

**30% stronger** on Skylake than on Haswell

## Comparison on Haswell – 2.5 GHz

Work	Impl.	CPU	Compiler	Key Generation		Shared Secret	
				Latency [cycles]	Throughput [ops/sec]	Latency [cycles]	Throughput [ops/sec]
[FHLD19]	(2 × 2)-way	i7-4770	Clang 5.0.2	43,700	57,208	121,000	20,661
	(2 × 2)-way	i7-4710HQ	Clang 10.0.0	41,938	59,575	121,499	20,563
[NS20]	(4 × 1)-way	i7-6500U	GCC 7.3.0	100,127	24,968	120,108	20,815
	(4 × 1)-way	i7-4710HQ	GCC 8.4.0	100,669	24,820	120,847	20,676
<b>This work</b>	(4 × 1)-way	i7-4710HQ	Clang 10.0.0	104,579	<b>95,568</b> 60.4%	329,455	<b>30,336</b> 45.7%

## Comparison on Skylake – 2.0 GHz

Work	Impl.	CPU	Compiler	Key Generation		Shared Secret	
				Latency [cycles]	Throughput [ops/sec]	Latency [cycles]	Throughput [ops/sec]
[FHLD19]	(2 × 2)-way	i7-6700K	Clang 5.0.2	34,500	57,971	99,400	20,150
	(2 × 2)-way	i5-6360U	Clang 10.0.0	35,629	55,955	95,129	20,939
[HEY20]	(4 × 1)-way	i9-7900X	GCC 5.4	n/a	n/a	98,484	20,308
	(4 × 1)-way	i5-6360U	GCC 8.4.0	n/a	n/a	116,595	16,656
[NS20]	(4 × 1)-way	i7-6500U	GCC 7.3.0	84,047	23,796	95,437	20,956
	(4 × 1)-way	i5-6360U	GCC 8.4.0	82,054	24,406	93,657	21,168
<b>This work</b>	(4 × 1)-way	i5-6360U	Clang 10.0.0	80,249	<b>99,363</b> 71.4%	246,636	<b>32,318</b> 52.7%

# Conclusion

- ⊙ AVX2 offers great potential to optimize ECC
- ⊙ The first to use AVX2 to maximize throughput
- ⊙ 1.5x ~ 1.7x throughput compared to the state of the art
- ⊙ Straightforward extension to AVX512

## Future Work

- ⦿ Support AVX512
- ⦿ Isogeny-based cryptography

Source code at  
<https://gitlab.uni.lu/APSIA/AVXECC>

Thank you for your attention!