



# NTT Multiplication for NTT-unfriendly Rings

New Speed Records for Saber and NTRU on Cortex-M4 and AVX2

Chi-Ming Marvin Chung Vincent Hwang Matthias J. Kannwischer Gregor Seiler Cheng-Jhih Shih Bo-Yin Yang

CHES 2021 - September 17, 2021

Number-theoretic Transforms

#### Number-theoretic Transforms

Suppose *n* is invertible.

Given an invertible  $\zeta \in R$  and a principle *n*-th root of unity  $\omega$  ( $\forall k, n \not| k \longrightarrow \sum_{i=0}^{n-1} \omega^{ik} = 0$ ), we have

NTT: 
$$R[x]/\langle x^n - \zeta^n \rangle \cong \prod_{i=0}^{n-1} R[x]/\langle x - \zeta \omega^i \rangle.$$

- Cyclic:  $\zeta^n = 1$
- Negacyclic:  $\zeta^n = -1$
- Cooley–Tukey:  $\zeta^n = (\omega^j)^n$ , the *n*-th power of a power of  $\omega$

As an isomorphism, we have

$$\begin{cases} \operatorname{NTT}(\boldsymbol{a}(x)\boldsymbol{b}(x)) = \operatorname{NTT}(\boldsymbol{a}(x))\operatorname{NTT}(\boldsymbol{b}(x)) & \iff \boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x))\operatorname{NTT}(\boldsymbol{b}(x))) \\ \operatorname{NTT}(\boldsymbol{a}(x) + \boldsymbol{b}(x)) = \operatorname{NTT}(\boldsymbol{a}(x)) + \operatorname{NTT}(\boldsymbol{b}(x)) & \iff \boldsymbol{a}(x) + \boldsymbol{b}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x)) + \operatorname{NTT}(\boldsymbol{b}(x))) \end{cases}$$



### **Fast Fourier Transforms**

- Let  $\omega$  be a principal  $n_0n_1$ -th root of unity.
  - Cooley–Tukey:

$$R[x]/\langle x^{n_0n_1} - 1 \rangle \cong \prod_{i=0}^{n_0-1} R[x]/\langle x^{n_1} - \omega^{n_1i} \rangle \cong \prod_{i=0}^{n_0-1} \prod_{j=0}^{n_1-1} R[x]/\langle x - \omega^{i+n_0j} \rangle$$

• Gentleman–Sande:

$$R[x]/\langle x^{n_0n_1} - 1 \rangle \cong \prod_{i=0}^{n_0-1} R[x]/\langle x^{n_1} - \omega^{n_1i} \rangle \cong \prod_{i=0}^{n_0-1} R[x]/\langle x^{n_1} - 1 \rangle \cong \prod_{i=0}^{n_0-1} \prod_{j=0}^{n_1-1} R[x]/\langle x - \omega^{n_0j} \rangle$$

• For  $q_0 \perp q_1$ , Good's trick:

$$R[x]/\langle x^{q_0q_1}-1\rangle \cong (R[z]/\langle z^{q_1}-1\rangle)[y]/\langle y^{q_0}-1\rangle$$



## **NTTs: Montgomery Multiplications**

- Cortex-M4: smull, smlal  $\implies$  32-bit arithmetic
- AVX2: vpmulhw ⇒ 16-bit arithmetic

Algorithm 1 32-bit Montgomery multiplication on Cortex-M4	Algorithm 2 16-bit Montgomery multiplication with AVX2			
<b>Input:</b> $(c0, c1) = (a, b)$	<b>Require:</b> $a \in [-2^{15}, 2^{15}), b \in [-$	$\frac{q-1}{2}, \frac{q-1}{2}], b' = bq^{-1} \mod 2^{16}$		
<b>Output:</b> $c0 \equiv ab2^{-32} \pmod{q}$	<b>Ensure:</b> $r \equiv ab2^{-16} \pmod{q}$			
	1: $t_1 \leftarrow \left\lfloor \frac{ab}{2^{16}} \right\rfloor$	▷ signed high product		
1: smull tmp0, c0, c0, c1	2: $t_0 \leftarrow ab' \mod 2^{16}$	▷ signed low product		
2: mul tmp1, tmp0, $(-q^{-1} \mod {}^{\pm}\mathrm{R})$	3: $t_0 \leftarrow \left\lfloor \frac{t_0 q}{2^{16}} \right\rfloor$	▷ signed high product		
3: smlal tmp0, c0, tmp1, $q$	4: $r \leftarrow (t_1 - t_0) \mod 2^{16}$			



Saber

#### Saber

- *q* = 8192
- $R_q = \mathbb{Z}_q[x]/\langle x^{256} + 1 \rangle$
- Parameter  $(l, \mu)$  varies for security levels:
  - Lightsaber:  $(I, \mu) = (2, 10)$
  - Saber:  $(I, \mu) = (3, 8)$
  - Firesaber:  $(I, \mu) = (4, 6)$
- Compute  $A^T \cdot s$  and  $A \cdot s'$  where
  - $A \in R_q^{l \times l}$
  - $s,s' \in R'_q$ , coefficients are in  $[-rac{\mu}{2},rac{\mu}{2}]$ , small



# How NTT-friendly/unfriendly Saber is?

- Polynomial modulus:  $x^{256} + 1$ , awesome!
- Coefficient ring:  $\mathbb{Z}_{8192}$ , unfriendly
  - Solution: choose ...
  - Hold on.

 $\implies$ 

Recall if NTT is defined correctly, then

 $\boldsymbol{a}(x)\boldsymbol{b}(x) + \boldsymbol{c}(x)\boldsymbol{d}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x))\operatorname{NTT}(\boldsymbol{b}(x)) + \operatorname{NTT}(\boldsymbol{c}(x))\operatorname{NTT}(\boldsymbol{d}(x)))$ 

$$A^T \cdot s = \operatorname{NTT}^{-1}(\operatorname{NTT}(A^T)\operatorname{NTT}(s))$$



# **NTTs for** $A^T \cdot s$

- We compute  $A^T \cdot s$  as if  $\mathbb{Z}$  is the coefficient ring
- Bounding the maximum value of the result:
  - Max:  $2 \cdot 256 \cdot \frac{8192}{2} \cdot \frac{\mu}{2} \cdot l = 2^{20} \cdot \mu \cdot l$
  - Cortex-M4: choose a 32-bit prime  $q' > 2^{20} \cdot \mu \cdot l$
  - AVX2: choose two 16-bit primes  $p_0, p_1$  with  $p_0p_1 > 2^{20} \cdot \mu \cdot I$
- Compute  $A^T \cdot s$  as
  - Cortex-M4:  $NTT^{-1}(NTT(A^T)NTT(s))$  in  $\mathbb{Z}_{q'}$
  - AVX2:  $\operatorname{NTT}^{-1}(\operatorname{NTT}(A^T)\operatorname{NTT}(s))$  in  $\mathbb{Z}_{p_0}$  and  $\mathbb{Z}_{p_1}$ , and CRT at the end
- $l + l^2$  NTTs: NTT(s) and NTT( $A^T$ )
- *l*<sup>2</sup> "point multiplications"
- $/ \operatorname{NTT}^{-1}: \operatorname{NTT}^{-1} (\operatorname{NTT}(A^T) \operatorname{NTT}(s))$



#### NTTs for Saber on Cortex-M4

- Incomplete NTTs giving 4-coefficient polynomials
- Long multiplication with accumulation: smlal
  - Let **p**<sub>i</sub>, **q**<sub>i</sub> be 4-coefficient polynomials
  - $\star$  is multiplication in (mod  $x^4 \zeta$ )
  - For  $\mathbf{h} = \sum_{i} \mathbf{p}_{i} \star \mathbf{q}_{i}$ ,  $[x^{0}]\mathbf{h} = \sum_{i} (\mathbf{p}_{i0}\mathbf{q}_{i0} + \zeta(\mathbf{p}_{i1}\mathbf{q}_{i3} + \mathbf{p}_{i2}\mathbf{q}_{i2} + \mathbf{p}_{i3}\mathbf{q}_{i1}))$
  - We can
    - montgomeryR for each of  $\mathbf{p}_{i0}\mathbf{q}_{i0} + \zeta(\mathbf{p}_{i1}\mathbf{q}_{i3} + \mathbf{p}_{i2}\mathbf{q}_{i2} + \mathbf{p}_{i3}\mathbf{q}_{i1})$ , or
    - montgomeryR for [x<sup>0</sup>]h



NTRU

### NTRU

- $\mathbb{Z}_q[x]/\langle x^n-1\rangle$
- Parameter (q, n) varies for security levels:
  - ntruhps2048509: (q, n) = (2048, 509)
  - ntruhps2048677: (q, n) = (2048, 677)
  - ntruhrss701: (q, n) = (8192, 701)
  - ntruhps4096821: (q, n) = (4096, 821)
  - ntruhps40961229: (q, n) = (4096, 1229)
  - ntruhrss1373: (q, n) = (16384, 1373)
- One of the multiplicands is ternary, i.e. coefficients are in  $\{\pm 1, 0\}$



# NTTs for NTRU on Cortex-M4

Parameter sets	NTT <sub>N</sub>	q'	Strategy
ntruhps4096821	$1728 = 9 \cdot 64 \cdot 3$	3365569	Mixed-radix (CT+GS)
ntruhrss701	$1536 = 512 \cdot 3$	5747201	Good's (CT+CT)
ntruhps2048677	$1536 = 512 \cdot 3$	1389569	Good's (CT+CT)
ntruhps2048509	$1024 = 256 \cdot 4$	1043969	Radix-2 (CT+GS)

In NTRU Prime, primes *p*, *q* give the field  $\mathbb{Z}_q[x]/\langle x^p - x - 1 \rangle$ . We compare (q, p) = (4591, 761) in NTRU Prime with (q, n) = (2048, 677) and (8192, 701) in NTRU.

- $\langle x^{\text{NTT}_{\mathbb{N}}} 1 \rangle \rightarrow \langle x^n 1 \rangle$  is faster than  $\langle x^{\text{NTT}_{\mathbb{N}}} 1 \rangle \rightarrow \langle x^p x 1 \rangle$ . Excluding head and tail cases,
  - $x^n \underline{1}$ : 1 add;  $q' > n \cdot q$
  - $x^p (x+1)$ : 2 adds;  $q' > n \cdot (2p-1)$
- $\mathbb{Z}_{q'} \to \mathbb{Z}_{\{2048, 8192\}}$  is faster than  $\mathbb{Z}_{q'} \to \mathbb{Z}_{4591}$ :
  - pkhbt before  $\mathbb{Z}_{q'} \to \mathbb{Z}_{\{2048, 8192\}}$ : 0.5 cycles on average with the and instruction
  - pkhbt after  $\mathbb{Z}_{q'} \to \mathbb{Z}_{4591}$ : 2 cycles on average with Barrett reduction
- On average, we save 2.5 cycles for each coefficient.
- Polynomials in NTRU are shorter than in NTRU Prime  $\implies > 2.5 * 701 = 1752.5$  cycles of reduction



Implementation Considerations with AVX2

## NTTs with AVX2

- Divided difference form of CRT:
  - Solve  $|u| < P/2 = \prod_{i=1}^{s} p_i/2$  from  $u \equiv u_i \pmod{p_i}$ ,  $i = 1 \dots s$ ,  $|u_i| < p_i/2$  Let  $m_i := (p_1 \cdots p_{i-1})^{-1} \mod \pm p_i$ .

$$\begin{cases} y_1 = u_1 \\ y_2 = y_1 + ((u_2 - y_1)m_2 \mod {}^{\pm}p_2) p_1 \\ y_3 = y_2 + ((u_3 - y_2)m_3 \mod {}^{\pm}p_3) p_1 p_2 \\ \vdots & \vdots \\ u = y_s = y_{s-1} + ((u_s - y_{s-1})m_s \mod {}^{\pm}p_s) p_1 \cdots p_{s-1} \end{cases}$$

- Range analysis
  - Compute worst-case intervals from input intervals
  - Precise operations and roots of unity

Results

#### Saber Results: MatrixVectorMul and InnerProd

MatrixVectorMul							
	Co	ortex-M4	Skylake (AVX2)				
	[BMKV20]	Ou	r Work	[BMKV20]	Our	Work	
<i>l</i> = 2	159k	66k	(- 58%)	7 002	5215	(-25%)	
<i>l</i> = 3	317k	125k	(-61%)	14 145	9 579	(-32%)	
<i>l</i> = 4	528k	205k	(-61%)	24 342	14 959	(-39%)	
InnerProduct							

**Table 1:** Cycles for MatrixVectorMul and InnerProd in Saber.

	Co	ortex-M4	Skylake (AVX2)		
	[BMKV20]	Our Work	[BMKV20]	Our Work	
<i>l</i> = 2	73k	41k (- 44%)	4 0 1 6	2125 (-47%)	
<i>l</i> = 3	99k	57k (- 42%)	5 977	2706 (-55%)	
<i>l</i> = 4	126k	73k (- 42%)	8 040	3278 (-60%)	



#### Saber Results: Full Scheme

		Cortex-M4		Skylake (AVX2)		
		[BMKV20]	Our Work	[BMKV20]	Our Work	
K Lightsaber E		466k	360k(-23%)	61 325	59831 (-2%)	
		653k	513k(-21%)	75 876	72473 (-4%)	
D	D	678k	498k(-27%)	70 228	64859 (-8%)	
I	K	853k	658 (-23%)	104 832	9971 (-5%)	
Saber	E	1103k	864 (-22%)	125 835	11844 (-6%)	
1	D	1127k	835 (-26%)	118 553	10726 (-10%)	
I	K	1 340k	1008k(-25%)	157 915	148729 (-6%)	
Firesaber	E	1642k	1255k(-24%)	184 322	171993 (-7%)	
I	D	1679k	1227k(-27%)	177 864	159950(-10%)	

**Table 2:** Clock cycles for Lightsaber, Saber, and Firesaber.



## **NTRU Results: Polynomial Multiplications**

Table 3: Clock cycles for big by small polynomial multiplication in NTRU.

	Cortex-M4			Skyl	ake (AVX	(2)
п	[KRS19]	Our Work		[ZCH <sup>+</sup> 19]	Our	Work
509	104k	101k	(- 3%)	6 6 4 3	8 540	(+29%)
677	175k	156k	(-11%)	11 103	10373	(-7%)
701	173k	156k	(-10%)	11 242	10 373	(-8%)
821	230k	199k	(- 13%)	15 507	13247	(-15%)



### **NTRU Results: Full Scheme**

 Table 4: Clock cycles for NTRU.

	Cortex-M4		Skylake (AVX2)	
	[KRS19]	Our Work	[ZCH <sup>+</sup> 19]	Our Work
K	79 682k	79660k(±0%)	208 653	218 887 (+5%)
ntruhps2048509 <b>E</b>	572k	564k(-1%)	71018	73176 (+3%)
D	545k	538k(-1%)	38 950	42953(+10%)
K	143 808k	143725k(±0%)	332 906	333278 (±0%)
ntruhps2048677 <b>E</b>	849k	821k(-3%)	96 293	95953 (±0%)
D	845k	818k(-3%)	59 169	58 406 (-1%)
K	154 477k	$154403k(\pm0\%)$	299 066	298505 (±0%)
ntruhrss701 E	403k	377k(-6%)	56616	56084 (-1%)
D	896k	871k(-3%)	62 503	61199 (-2%)
K	208 953k	207495k(-1%)	458 614	451664 (-2%)
ntruhps4096821 <b>E</b>	1069k	1027k(-4%)	114 986	113 935 (-1%)
D	1075k	1030k(-4%)	74 182	70917 (-4%)



#### LAC Results: Polynomial Multiplications

- q = 251 for LAC-{128, 192, 256}-v3a
- $\mathbb{Z}_q[x]/\langle x^n+1\rangle$
- Parameter *n* varies for security levels:
  - LAC-128-v3a: *n* = 512
  - LAC-192-v3a and LAC-256-v3a: n = 1024
- One of the multiplicands is ternary, i.e. coefficients are in  $\{\pm 1, 0\}$

	Cortex-M4			Skyl	ake (AV>	<2)
	[LLZ <sup>+</sup> 18]	Our Work		[LLZ <sup>+</sup> 18]	Our	Work
LAC-128-v3a	638k	65k	(-90%)	14 691	4 552	(-69%)
LAC-192-v3a	1 274k	131k	(-90%)	73 955	10119	(-86%)
LAC-256-v3a	1701k	132k	(-92%)	73 955	10119	(-86%)

Table 5: LAC polynomial multiplication clock cycles on Cortex-M4 and Skylake



## LAC Results: Full Scheme

		Cortex-M4		Sky	lake (AVX2)
		[LLZ <sup>+</sup> 18]	Our Work	[LLZ <sup>+</sup> 18]	Our Work
	κ	850k	282k(-67%)	53 000	42167(-20%)
LAC-128-v3a	E	1 430k	450k(-69%)	76418	59252(-22%)
	D	1 960k	565k(-71%)	86 209	55 880 (-35%)
	κ	1 507k	373k(-75%)	96 270	41713(-57%)
LAC-192-v3a	E	2 427k	610k(-75%)	128 342	67732(-47%)
	D	3 329k	824k(-75%)	189 660	74393(-61%)
LAC-256-v3a	κ	2 020k	459k(-77%)	143 568	76917(-46%)
	E	3 633k	748k(-79%)	202 346	106836(-47%)
		5 327k	1111k(-79%)	262 901	104 897 (-60%)

 Table 6:
 Performance results in clock cycles for LAC



# Conclusion

#### Saber

- Coefficient ring: NTT-unfriendly
- Polynomial modulus: NTT-friendly
- MatrixVectorMul: NTT-friendly
- NTRU
  - Coefficient rings: NTT-unfriendly
  - The degrees of polynomials: large enough for NTTs
- LAC
  - Coefficient rings: NTT-unfriendly
  - Polynomial modulus: NTT-friendly
  - The degrees of polynomials: large enough for NTTs
- Compute the result as if  $\ensuremath{\mathbb{Z}}$  is the coefficient ring
  - Cortex-M4: 32-bit prime
  - AVX2: two 16-bit primes



# Works Worth Noting

- We optimized Saber on Cortex-M4 only for speed. We thank Michiel van Beirendonck for integrating stack optimizations. See commit 992f0f226503d43b6d33278ecb60a9168ed8d787 in pqm4.
- For even more stack optimized Saber on Cortex-M4 and more analysis on NTTs:
  - "Multi-moduli NTTs for Saber on Cortex-M3 and Cortex-M4"
  - https://eprint.iacr.org/2021/995
- For NTTs with 64-bit Armv8-A:
  - "Neon NTT: Faster Dilithium, Kyber, and Saber on Cortex-A72 and Apple M1"
  - https://eprint.iacr.org/2021/986
  - Barrett multiplication: 3-instruction single-width modular multiplication
  - Asymmetric multiplication: multiplying in  $R[x]/\langle x^{\{2,4\}} \zeta \rangle$  essentially as in  $R[x]/\langle x^{\{2,4\}} 1 \rangle$  without requiring  $\zeta = \omega^{\{2,4\}}$ , applicable whenever incomplete NTTs are re-used, e.g. Kyber and Saber





### Reference i

- Jose Maria Bermudo Mera, Angshuman Karmakar, and Ingrid Verbauwhede.
  Time-memory trade-off in toom-cook multiplication: an application to module-lattice based cryptography.
  IACR Transactions on Cryptographic Hardware and Embedded Systems, 2020(2):222–244, Mar. 2020.
- Matthias J. Kannwischer, Joost Rijneveld, and Peter Schwabe.
   Faster multiplication in Z<sub>2<sup>m</sup></sub>[x] on cortex-m4 to speed up NIST PQC candidates.
   In Applied Cryptography and Network Security, pages 281–301, 2019.
- Xianhui Lu, Yamin Liu, Zhenfei Zhang, Dingding Jia, Haiyang Xue, Jingnan He, and Bao Li.
   LAC: practical ring-lwe based public-key encryption with byte-level modulus.
   IACR Cryptol. ePrint Arch., 2018.
   https://eprint.iacr.org/2018/1009.
- Zhenfei Zhang, Cong Chen, Jeffrey Hoffstein, William Whyte, John M. Schanck, Andreas Hulsing, Joost Rijneveld, Peter Schwabe, and Oussama Danba.

## NTRU.

Technical report, National Institute of Standards and Technology, 2019.

available at https://csrc.nist.gov/projects/post-quantum-cryptography/round-2-submissions.

