

Multi-moduli NTTs for Saber on Cortex-M3 and Cortex-M4

Amin Abdulrahman Jiun-Peng Chen Yu-Jia Chen Vincent Hwang Matthias J. Kannwischer Bo-Yin Yang

CHES 2022, Leuven, Belgium

Organization of This Talk

Contributions

Backgrounds

Time–Memory Tradeoffs

Polynomial Multiplications on Cortex-M4

MatrixVectorMul

First-Order Masked MatrixVectorMul and InnerProd

Saber on Cortex-M3

Performance



Contributions

Contributions

- Time-memory tradeoffs for NTT-based approaches
 - Multiplying two polys.: 16-bit NTT + 32-bit NTT
 - MatrixVectorMul: cache NTT(s) or not, acc. in NTT domain or not
- First-order masked MatrixVectorMul and InnerProd
 - NTT(A) as before
 - Compute the shares of NTT(s) = NTT(s₀) + NTT(s₁)
- Cortex-M3
 - Both 16-bit and 32-bit NTTs are significantly slower because of
 - The absence of DSP instructions
 - The non-applicability of long multiplications to secrete data
 - Computing two 16-bit NTTs becomes faster



Results

- Cortex-M4
 - Cycles: NTT (speed) << NTT (stack) \approx non-NTT (speed: TMVP, Toom–Cook) << non-NTT (stack: Karatsuba)
 - Stack: NTT (stack) \approx non-NTT (stack) < non-NTT (speed) < NTT (speed)
- Cortex-M3
 - Cycles: 16-bit NTT (speed) < 32-bit NTT (speed) < 16-bit NTT (stack) < non-NTT (speed, Toom–Cook)
 - Stack: 16-bit NTT (stack) < 32-bit NTT (speed) \approx 16-bit NTT (speed) < non-NTT (speed, Toom–Cook)



Backgrounds

Saber

- $R_q = \mathbb{Z}_{8192}[x]/\langle x^{256}+1 \rangle$
- Parameters (I, μ) varies from security levels (other parameters omitted in this talk).
 - LightSaber: (/, μ) = (2,10)
 - Saber: $(I, \mu) = (3, 8)$
 - FireSaber : (*l*, μ) = (4, 6)
- $A \in R_q^{l \times l}$, s, s' $\in R_q^l$.
 - Key generation: $A^T s$
 - Encryption: As'



NTT-Based MatrixVectorMul for Saber

- Find an NTT-friendly modulus q' such that $A^T s$ in \mathbb{Z} is the same as in $\mathbb{Z}_{q'}$
 - NTT-firendly: next slide
 - Signed arithmetic: choose $q' > 2 \cdot \frac{8192}{2} \cdot \frac{\mu}{2} \cdot I$
- Compute $A^T s = \operatorname{NTT}^{-1} (\operatorname{NTT}(A^T) \cdot \operatorname{NTT}(s))$
 - *l*² + / NTTs
 - / NTT⁻¹s
 - l^2 base multiplications



Number-Theoretic Transforms i

- Ring *R*, invertible $\zeta \in R$
- $n \perp \operatorname{char}(R)$, principal *n*-th root of unity ω_n ($\forall 1 \leq i < n, \sum_{j=0}^{n-1} \omega_n^{ij} = 0$). Equivalently, for $R = \mathbb{Z}_q$ with prime factorization $q = \prod_{i=0}^{l-1} p_i^{d_i}$, $n | \mathbf{0}(q) \coloneqq \gcd(p_i 1)_{0 \leq i < l}$ [AB74].

- $R[x]/\langle x^n \zeta^n \rangle \cong \prod_{i=0}^{n-1} R[x]/\langle x \zeta \omega_n^i \rangle : \mathbf{a}(x) \mapsto \mathbf{a}(\zeta \omega_n^i)_i$
- Cooley–Tukey FFT: $O(n \log n)$ for $n = 2^k$

Number-Theoretic Transforms ii

- $R=\mathbb{Z}_{q_0q_1}$, $R_0=\mathbb{Z}_{q_0}$, $R_1=\mathbb{Z}_{q_1}$, $q_0ot q_1$
- $(\zeta_0, \zeta_1) = (\zeta \mod q_0, \zeta \mod q_1), (\omega_{0:n}, \omega_{1:n}) = (\omega_n \mod q_0, \omega_n \mod q_1)$
- NTT := $\boldsymbol{a}(x) \mapsto \boldsymbol{a}(\zeta \omega_n^i)_i$, NTT₀ := $\boldsymbol{a}(x) \mapsto \boldsymbol{a}(\zeta_0 \omega_{0:n}^{i_0})_{i_0}$, NTT₁ := $\boldsymbol{a}(x) \mapsto \boldsymbol{a}(\zeta_1 \omega_{1:n}^{i_1})_{i_1}$



HES 2022, 9/20 Vincent Hv



Time–Memory Tradeoffs

Memory for Polynomials

- Memory for buffers.
- Memory for public polynomials (on-the-fly generation).
- We ignore the memory for secrete polynomials.



Memory Usage of 32-bit NTT and 16-bit NTTs Approaches

Each line segment = 4096 bits (a size-256 poly with 16-bit coeffs.)

Figure 1: 32-bit.



CHES 2022, 9/20 Vincent



Observations of 16-bit and 32-bit NTTs

For a Cortex-M4, one 32-bit NTT is much faster than two 16-bit NTTs.

- 1. Start with 16-bit NTTs
- 2. Identify at which point that inevitably, corresponding elements in $\mathbb{Z}_{q_0}, \mathbb{Z}_{q_1}$ are *both* in memory
- 3. Replace operations in $\mathbb{Z}_{q_0}, \mathbb{Z}_{q_1}$ with $\mathbb{Z}_{q_0q_1}$ for these elements
- A speed optimization for the heavily stack-optimized implementation
 - Same memory usage as purely 16-bit approach
 - Faster than 16-bit approach

Performance of NTT-Based Polynomial Multiplications on Cortex-M4

 Table 1: NTT-related functions on Cortex-M4. Numbers of the last two columns are extracted for comparisons.

32-bit	16-bit + 16-bit	32-bit	16-bit
5 853	4 374+ 4 822	5853	4822
7 1 37	-	7137	4817
_	3731 + 2965	4186	2965
_	0+1171	-	-
_	2 435	-	2 4 3 5
	32 488	23 0 29	37 287
	1 536	2 0 4 8	1 536
	32-bit 5 853 7 137 – – –	32-bit 16-bit + 16-bit 5853 4374+ 4822 7137 - - 3731 + 2965 - 0 + 1171 - 2435 32488 1536	32-bit 16-bit + 16-bit 32-bit 5853 4374+ 4822 5853 7137 - 7137 - 3731 + 2965 4186 - 0 + 1171 - - 2435 - 32488 23029 1536 2048

Strategies for MatrixVectorMul

Figure 3: Strategies for MatrixVectorMul.

	Cache NTT(s)	Re-compute NTT(<i>s</i>)
Acc. in NTT domain	А	С
Acc. in $\mathbb{Z}_{8192}[x]$	В	D

- Key generation, A^Ts : strategies A, B, D
- Encryption, As': strategies A, C, D

First-Order Masked MatrixVectorMul and InnerProd

First-Order Masked MatrixVectorMul and InnerProd

- Split $s' = s'_0 + s'_1$ (first-order)
- Compute (As'_0, As'_1)
- $(As'_0, As'_1) = (NTT^{-1}(NTT(A) \cdot NTT(s'_0)), NTT^{-1}(NTT(A) \cdot NTT(s'_1)))$
 - $l^2 + 2l$ NTTs
 - 2/NTT⁻¹s
- Coefficient rings of s'_0, s'_1 : \mathbb{Z}_{8192} instead of $\left\{-\frac{\mu}{2}, \dots, \frac{\mu}{2}\right\}$
 - Compute with one 32-bit NTT and one 16-bit NTT
- In total:
 - *l*² + 2/ 32-bit NTTs
 - *l*² + 2/ 16-bit NTTs
 - 2/ 32-bit NTT⁻¹s
 - 2/16-bit NTT⁻¹s

Saber on Cortex-M3

Differences Between Cortex-M3 and Cortex-M4

- No floating-point registers
- No DSP extension (s{mul, mla}{b, t}{b, t}, smlad{, x}, {u, s}{add, sub}{8, 16})
 - 16-bit NTTs are much slower
- {u, s}{mul, mla}l takes input-dependent cycles
 - NTT_leak: 32-bit NTTs are variable time (for public data)
 - Constant-time NTTs: Emulate 32-bit NTTs with mul, mla, ... [GKS21] (much slower)
- Question: which is better?
 - Cortex-M4: one 32-bit NTT is faster than two 16-bit NTTs
 - Cortex-M3: two 16-bit NTTs vs one 32-bit NTT



Saber on Cortex-M3 i

- 16-bit NTTs only
 - $As' = NTT^{-1}(NTT(A) \cdot NTT(s'))$ where NTT/NTT^{-1} is a pair of 16-bit NTT/iNTTs
- 32-bit NTTs only
 - $As' = NTT^{-1}(NTT_{leak}(A) \cdot NTT(s'))$ where NTT is the constant-time NTT.

15/19

Saber on Cortex-M3 ii

Table 2: NTT-related functions on Cortex-M3.

	2 imes 16-bit	32-bit
NTT	16774	31 056
NTT_leak	-	19 363
NTT^{-1}	19079	37 394
base_mul	11 933	8 532
mod p _i	-	_
CRT	4 642	_
poly_mul	69 202	96 345

CHES 2022, 9/20



Performance

Cortex-M4 Results i

Table 3: Unprotected Saber on Cortex-M4.

			LightSaber		Saber		FireSaber	
			сс	stack	сс	stack	сс	stack
		к	612k	3 564	1 230k	4 348	2046k	5116
	[MKV20]	Е	880k	3148	1616k	3 412	2 538k	3 668
	(stack)	D	976k	3 164	1 759k	3 4 2 0	2740k	3 684
		κ	360k	14 604	658k	23 284	1008k	37 116
N.1.4	[CHK+21]	Е	513k	16252	864k	32 620	1 255k	40 484
	(speed)	D	498k	16996	835k	33 824	1 227k	41 964
1014	This work	К	353k	5764	644k	6788	990k	7812
	32-bit	Е	487k	6444	826k	7 468	1 208k	8 484
	(speed, A)	D	456k	6440	777k	7 484	1 152k	8 500
	This work	κ	423k	3 4 2 8	819k	3 940	1 315k	4 4 5 2
	hybrid	Е	597k	3 204	1063k	3 3 3 2	1617k	3 468
	(stack, D)	D	583k	3 2 2 0	1039k	3 348	1594k	3 4 8 4

CHES 2022, 9/20



Cortex-M4 Results ii

Table 4: Masked Saber (I = 3) on the Cortex-M4.

Table 5: Masking cycles/stack overhead.

	Decapsulation		
	cc stack		
[VBDK+20]	2 833k	11656	
This work (speed, A)	2 385k	16 140	
This work (C)	2615k	10 476	
This work (stack, D)	2846k	8432	

	unma	sked A	unma	sked D
	cc stack		сс	stack
masked A	3.07	2.16	2.30	4.82
masked C	3.37	1.40	2.52	3.13
masked D	3.66	1.13	2.74	2.52

Cortex-M3 Results

Table 6: Unprotected Saber on Cortex-M3.

			LightSaber		Saber		FireSaber	
			сс	stack	сс	stack	сс	stack
	pqm3	к	710k	9652	1 328k	13 252	2171k	20 1 16
	Toom	Е	967k	11372	1738k	15516	2 688k	22 964
	(speed)	D	1081k	12116	1 902k	16612	2 933k	24 444
	This work	κ	540k	5756	939k	6 788	1439k	7 812
М3	16-bit	Е	715k	6436	1194k	7 468	1751k	8 4 9 2
	(speed, A)	D	749k	6436	1 237k	7 468	1811k	8 4 9 2
	This work	κ	632k	3 4 2 0	1 253k	3 9 3 2	1 955k	4 4 4 4
	16-bit	Е	887k	3 2 0 4	1614k	3 332	2 427k	3 460
	(stack, D)	D	923k	3 2 0 4	1 657k	3 3 3 2	2 487k	3 460
	This work	К	594k	5732	1057k	6756	1 553k	7 788
	32-bit	Е	800k	6412	1 330k	7 444	1 883k	8 468
	(speed, A)	D	877k	6420	1 429k	7 452	2016k	8 476

CHES 2022, 9/20



Reference i

🔋 RC Agarwal and C Burrus.

Fast convolution using Fermat number transforms with applications to digital filtering.

IEEE Transactions on Acoustics, Speech, and Signal Processing, 22(2):87–97, 1974.

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

NTT Multiplication for NTT-unfriendly Rings New Speed Records for Saber and NTRU on Cortex-M4 and AVX2.

IACR Transactions on Cryptographic Hardware and Embedded Systems, 2021(2):159–188, 2021.

https://tches.iacr.org/index.php/TCHES/article/view/8791.



Reference ii

Denisa O. C. Greconici, Matthias J. Kannwischer, and Daan Sprenkels.
 Compact Dilithium Implementations on Cortex-M3 and Cortex-M4.
 IACR Transactions on Cryptographic Hardware and Embedded Systems, 2021(1):1–24, 2021.

https://tches.iacr.org/index.php/TCHES/article/view/8725.

 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, 2020.
 https://tches.iacr.org/index.php/TCHES/article/view/8550.



Reference iii

Michiel Van Beirendonck, Jan-Pieter D'Anvers, Angshuman Karmakar, Josep Balasch, and Ingrid Verbauwhede.

A side-channel resistant implementation of SABER.

IACR Cryptol. ePrint Arch., 2020:733, 2020.

https://eprint.iacr.org/2020/733.



Cooley-Tukey FFT i

$$R[x]/\langle x^{n_0n_1} - \zeta^{n_0n_1} \rangle \cong \prod_{i_0=0}^{n_0-1} R[x]/\langle x^{n_1} - \zeta^{n_1} \omega_n^{i_0n_1} \rangle$$
$$\cong \prod_{i_0=0}^{n_0-1} \prod_{i_1=0}^{n_1-1} R[x]/\langle x - \zeta \omega_n^{i_0+i_1n_0} \rangle$$



CHES 2022, 9/20 Vincent Hwar

Cooley-Tukey FFT ii

$$R[x] / \langle x^{2^{k}} - \zeta^{2^{k}} \rangle \cong \prod_{i_{0}=0}^{1} R[x] / \langle x^{2^{k-1}} - \zeta^{2^{k-1}} \omega_{2}^{i_{0}} \rangle$$
$$\cong \prod_{i_{0}, i_{1}=0}^{1} R[x] / \langle x^{2^{k-2}} - \zeta^{2^{k-2}} \omega_{4}^{i_{0}+2i_{1}} \rangle$$
$$\cong \prod_{i_{0}, \dots, i_{k-1}=0}^{1} R[x] / \langle x - \zeta \omega_{2^{k}}^{\sum_{j=0}^{k-1} 2^{j} i_{j}} \rangle$$

- *k* isomorphisms
- Each isomorphism takes $O(2^k)$ time $\implies O(k2^k)$ time (or $O(n \lg n)$ where $n = 2^k$)



- Each line segment = 4096 bits, 16384 bits in total.
- A size-256 poly. of 32-bit coeffs. is stored in two segments.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}(\boldsymbol{b}(x)))$
- Expand *a*(*x*), *b*(*x*) to 32-bit first





- Each line segment = 4096 bits, 16384 bits in total.
- A size-256 poly. of 32-bit coeffs. is stored in two segments.
- Compute a(x)b(x) = NTT⁻¹(NTT(a(x)) · NTT(b(x)))
- Expand *a*(*x*), *b*(*x*) to 32-bit first



NTT(a(x))



- Each line segment = 4096 bits, 16384 bits in total.
- A size-256 poly. of 32-bit coeffs. is stored in two segments.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}(\boldsymbol{b}(x)))$
- Expand a(x), b(x) to 32-bit first



 $NTT(\mathbf{a}(x))$

 $NTT(\mathbf{b}(x))$



CHES 2022, 9/20 Vincent Hwa

- Each line segment = 4096 bits, 16384 bits in total.
- A size-256 poly. of 32-bit coeffs. is stored in two segments.
- Compute a(x)b(x) = NTT⁻¹(NTT(a(x)) · NTT(b(x)))
- Expand *a*(*x*), *b*(*x*) to 32-bit first



 $NTT(\boldsymbol{a}(x)) \cdot NTT(\boldsymbol{b}(x))$

CHES 2022, 9/20 Vincent Hwa



- Each line segment = 4096 bits, 16384 bits in total.
- A size-256 poly. of 32-bit coeffs. is stored in two segments.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}(\boldsymbol{b}(x)))$
- Expand *a*(*x*), *b*(*x*) to 32-bit first



 $\operatorname{NTT}^{-1}(\operatorname{NTT}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}(\boldsymbol{b}(x)))$

CHES 2022, 9/20 Vincent Hw



- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \mathsf{CRT}\left(\mathsf{NTT}_{0}^{-1}(\mathsf{NTT}_{0}(\boldsymbol{a}(x)) \cdot \mathsf{NTT}_{0}(\boldsymbol{b}(x))), \mathsf{NTT}_{1}^{-1}(\mathsf{NTT}_{1}(\boldsymbol{a}(x)) \cdot \mathsf{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$

CHES 2022, 9/20 Vincen



- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{CRT}\left(\operatorname{NTT}_{0}^{-1}(\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x))), \operatorname{NTT}_{1}^{-1}(\operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$



- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{CRT}\left(\operatorname{NTT}_{0}^{-1}(\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x))), \operatorname{NTT}_{1}^{-1}(\operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$



 $NTT_0(\boldsymbol{a}(x)) \cdot NTT_0(\boldsymbol{b}(x)) \qquad NTT_1(\boldsymbol{a}(x))$

CHES 2022, 9/20 Vincent Hw

- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{CRT}\left(\operatorname{NTT}_{0}^{-1}(\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x))), \operatorname{NTT}_{1}^{-1}(\operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$



- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{CRT}\left(\operatorname{NTT}_{0}^{-1}(\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x))), \operatorname{NTT}_{1}^{-1}(\operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$



 $\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x)) \operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x))$

CHES 2022, 9/20 Vincent Hwa



- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute $\boldsymbol{a}(x)\boldsymbol{b}(x) = \operatorname{CRT}\left(\operatorname{NTT}_{0}^{-1}(\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x))), \operatorname{NTT}_{1}^{-1}(\operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x)))\right)$
- Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$





- Each line segment = 4096 bits, 12288 bits in total.
- A size-256 poly. of 16-bit coeffs. is stored in a segment.
- Compute

 $\boldsymbol{a}(x)\boldsymbol{b}(x) = \mathsf{CRT}\left(\mathsf{NTT}_0^{-1}(\mathsf{NTT}_0(\boldsymbol{a}(x)) \cdot \mathsf{NTT}_0(\boldsymbol{b}(x))), \mathsf{NTT}_1^{-1}(\mathsf{NTT}_1(\boldsymbol{a}(x)) \cdot \mathsf{NTT}_1(\boldsymbol{b}(x)))\right)$

• Notice $(\boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_0, \boldsymbol{a}(x)\boldsymbol{b}(x) \mod q_1) = (\operatorname{NTT}_0^{-1}(\operatorname{NTT}_0(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_0(\boldsymbol{b}(x))), \operatorname{NTT}_1^{-1}(\operatorname{NTT}_1(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_1(\boldsymbol{b}(x))))$



 $a(x)b(x) \mod q_0q_1$







NTT(a(x))



CHES 2022, 9/20 Vincent Hwa





CHES 2022, 9/20



CHES 2022, 9/20 Vincent Hwa





CHES 2022, 9/20 Vincent Hwa

lwang Institute





CHES 2022, 9/20 Vincent Hwang





CHES 2022, 9/20 Vincent Hwa

 $\operatorname{NTT}_{0}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{0}(\boldsymbol{b}(x)) \operatorname{NTT}_{1}(\boldsymbol{a}(x)) \cdot \operatorname{NTT}_{1}(\boldsymbol{b}(x))$



2022, 9/20 Vincent Hwang



 $NTT(\boldsymbol{a}(x)) \cdot NTT(\boldsymbol{b}(x))$



CHES 2022, 9/20 Vincent Hwan

 $a(x)b(x) \mod q_0q_1$



CHES 2022, 9/20 Vincent Hwan

Combining 16-bit and 32-bit NTTs on Cortex-M3

- As' = NTT⁻¹((a → (a mod q₀, a mod q₁) ∘ NTT_leak)(A) · NTT(s')) where NTT/NTT⁻¹ is a pair of 16-bit NTT/iNTTs
- Doesn't worth it as $a\mapsto (a \mod q_0, a \mod q_1)pprox ext{NTT} ext{NTT}_ ext{leak}$

