

# PQC-AMX: Accelerating Saber and FrodoKEM on the Apple M1 and M3 SoCs

Décio Luiz Gazzoni Filho<sup>\*†</sup>, Guilherme Brandão<sup>‡</sup>, Gora Adj<sup>§</sup>,  
Arwa Alblooshi<sup>§</sup>, Isaac A. Canales-Martínez<sup>§</sup>, Jorge Chávez-Saab<sup>§</sup> and Julio López<sup>\*</sup>

<sup>\*</sup> Instituto de Computação, Universidade Estadual de Campinas (UNICAMP), Campinas, Brazil  
Email: {decio.gazzoni,jlopez}@ic.unicamp.br

<sup>†</sup> Dept. of Electrical Engineering, State University of Londrina, Londrina, Brazil. Email: dgazzoni@uel.br

<sup>‡</sup> Independent Researcher, Londrina, Brazil. Email: brandaogbs@gmail.com

<sup>§</sup> Cryptography Research Centre, Technology Innovation Institute, Abu Dhabi, UAE  
Email: {gora.adj,arwa.alblooshi,isaac.canales,jorge.saab}@tii.ae

**Abstract**—As CPU performance cannot keep up with the dramatic growth of the past few decades, CPU architects turn to domain-specific architectures to accelerate certain tasks. A recent trend is the introduction of matrix-multiplication accelerators to CPUs by manufacturers such as IBM, Intel and ARM, some of them yet to launch commercially. Apple’s systems-on-chip (SoCs) for its mobile phones, tablets and personal computers include a proprietary, undocumented CPU-coupled matrix multiplication coprocessor called AMX. We leverage AMX to accelerate the post-quantum lattice-based cryptosystems Saber and FrodoKEM, and benchmark their performance on Apple M1 and M3 SoCs. We propose a variant of the Toeplitz Matrix-Vector Product algorithm for polynomial multiplication, which sets new speed records for Saber using AMX, improving up to 20% for the main KEM operations, and 152% for matrix-vector multiplication of polynomials, over the current state-of-the-art. We also set new FrodoKEM speed records using AMX, gaining up to 21% for the main KEM operations and 124% for matrix multiplication (with further improvements for  $4\times$ -batching), over our optimized NEON implementation, also introduced here, which already improves upon the previous state-of-the-art for ARMv8 CPUs.

**Index Terms**—Post-quantum cryptography, AMX, ARM, NEON, FrodoKEM, Saber

## I. INTRODUCTION

Quantum computers pose a threat to cryptosystems whose security relies on the presumed hardness of computational problems such as integer factoring and discrete logarithms. Post-quantum cryptography (PQC) refers to cryptosystems that remain secure against attacks employing quantum and classical computers. In 2017, the National Institute of Standards and Technology (NIST) called for a PQC standardization process; 3 out of 4 selected candidates for standardization are lattice-based. Saber [1] and FrodoKEM [2] are lattice-based Key Encapsulation Mechanisms that reached round 3 in the standardization process; the latter is recommended by the German BSI [3] and is being considered for standardization by ISO [4].

The performance bottlenecks of lattice-based cryptography usually lie in polynomial/matrix multiplication and symmetric-cryptography operations, prompting extensive research efforts to enhance their efficiency. Various manufacturers have developed high-performance AI accelerators, such as NVIDIA’s tensor cores [5], Intel’s Advanced Matrix Extensions (AMX)

[6] and ARMv9-A’s Scalable Matrix Extensions (SME) [7], to cater to the high demands of AI applications. Apple’s AMX (unrelated to Intel’s) is an undocumented coprocessor found in its SoCs, starting with the 2019’s A13 [8, §7.6], which we apply to Saber and FrodoKEM in this work.

### A. Related works

Many studies target GPU cores, achieving high throughput via huge batching levels, but compromising latency. Gazzoni Filho et al. [9] presented the first cryptographic implementation on a CPU-linked matrix multiplication accelerator, setting new NTRU speed records using Apple’s AMX coprocessor on M1/M3 SoCs, beating state-of-the-art NEON implementations with low latency, no batching and running in constant time.

Becker et al. [10] set the current speed record for Saber [1] on ARMv8-A using  $\mathcal{O}(n \log n)$  Number Theoretical Transform (NTT) methods combined with a novel “Barrett multiplication” algorithm for modular multiplication, achieving a speedup of 56% over the previous state-of-the-art on the Apple M1. We also remark the work in [11], which introduced innovative Toeplitz Matrix-Vector Product (TMVP) formulas, with the “four-way” formula standing out as the best non-NTT-like multiplication algorithm for Saber’s ring on Cortex-M4.

The state-of-the-art implementation for FrodoKEM [2] is the work of Bos et al. [12], which improves matrix multiplication through a row-wise blocking and packing approach, and also proved that Strassen’s algorithm improves throughput for use cases with high batching levels. An ARMv8 implementation using NEON was presented shortly after in [13], claiming a speedup of  $10.22\times$  at the protocol level. However, they do not acknowledge the improvements of [12] which, while not ARM-specific, appear to be superior on M1 and M3.<sup>1</sup>

### B. Our contributions

We present an AMX implementation of Saber, adapting the techniques from [9] and modifying the TMVP method of [14] to benefit from batched products. Note that the improved

<sup>1</sup>The implementation of [13] is not publicly available and the authors could not be reached for clarification. Our benchmarks of §V show speedups for [12] that exceed those reported by [13] for their implementation.

TMVP formulas from [11] require more matrix multiplications but of a smaller size, but for AMX even the original formulas underutilize the throughput and there is no clear benefit from smaller matrices. Therefore, we focus only on the original TMVP formulas of [14] and increase throughput by batching. This sets new speed records on Apple M1 and M3, with speedups of up to 13% at the protocol level and 151% for polynomial operations.

For FrodoKEM, we introduce a new NEON implementation of our own to use as a baseline, which already sets new speed records on the M1/M3. We then present our AMX implementation, which improves further on our NEON record. Both implementations explore possible matrix multiplication strategies and use a novel technique for generating FrodoKEM-AES’s  $\mathbf{A}$  matrix. We make an innovative use of AMX’s unique `genlut` instruction to perform Gaussian sampling, improving it by up to 418% versus a NEON implementation. This might be of particular interest for other applications. Compared to the state of the art, we improve on the M1 and M3 by up to 21% at the protocol level and 124% for matrix multiplication. Then, we develop 4×-batched NEON and AMX implementations, showing that AMX is significantly faster than NEON, by up to 91% at the protocol level and 708% for matrix multiplication.

Our code is available under an open source CC0 license at <https://github.com/dgazzoni/PQC-AMX>.

## II. PRELIMINARIES

A public-key encryption scheme (PKE) is a tuple of algorithms ( $\text{KeyGen}, \text{Enc}, \text{Dec}$ ).  $\text{KeyGen}$  generates a public key  $pk$  and a secret key  $sk$ .  $\text{Enc}$  outputs a ciphertext  $c$  given  $pk$  and a message  $m$ .  $\text{Dec}$  outputs a message  $m'$  from  $sk$  and  $c$ . A key encapsulation mechanism (KEM) is a tuple of algorithms ( $\text{KeyGen}, \text{Encaps}, \text{Decaps}$ ).  $\text{KeyGen}$  generates a public key  $pk$  and a secret key  $sk$ .  $\text{Encaps}$  outputs a shared key  $ss$  and a ciphertext  $c$  given  $pk$ .  $\text{Decaps}$  outputs a shared key  $ss'$  from  $sk$  and  $c$ . We present next KEMs obtained from PKEs via a variant of the Fujisaki-Okamoto transform; we only show PKE algorithms, which are the target of our optimizations.

Bold lower case denotes vectors and bold upper case denotes matrices. We write  $\mathbf{v}[i : j : k]$  for a matrix/vector slice of coefficients  $i, i+j, i+2j, \dots, i+k$ ;  $j = 1$  if omitted. Sampling from a uniform distribution over a set  $S$  is denoted  $x \leftarrow \mathcal{U}(S)$ .

### A. Saber

Saber [1] is a lattice-based KEM relying on the hardness of Module Learning With Rounding. Its NIST submission specifies the parameter sets below for security levels 1, 3, and 5.

Parameter set	Sec. level	$l$	$n$	$q$	$p$	$T$	$\mu$
LightSaber	1	2	256	$2^{13}$	$2^{10}$	$2^3$	10
Saber	3	3	256	$2^{13}$	$2^{10}$	$2^4$	8
FireSaber	5	4	256	$2^{13}$	$2^{10}$	$2^6$	6

Saber works over  $R_q := \mathbb{Z}_q[X]/(X^n + 1)$  and employs the binomial distribution centered at  $\mu$ , denoted  $\beta_\mu$ , hash functions  $\mathcal{F}, \mathcal{G}, \mathcal{H}$ , and a function `gen` to generate a pseudorandom matrix from a seed.  $q, p$  and  $T$  are powers of 2 with exponents

$\epsilon_q, \epsilon_p, \epsilon_T \in \mathbb{Z}$ , resp. Let  $\mathbf{s} \leftarrow \beta_\mu(R_q^l; r)$  denote sampling each coordinate of a vector  $\mathbf{s} \in R_q^l$  pseudorandomly from the distribution  $\beta_\mu(R_q)$  with seed  $r$ . Algorithms II.1, II.3 and II.2 are a verbatim reproduction of Saber’s PKE specification.

Algorithm	II.1	Algorithm	II.3
Saber.PKE.KeyGen()		Saber.PKE.Enc( $pk, m; r$ )	
<b>Input:</b> None <b>Output:</b> Key pair ( $pk, sk$ )		<b>Input:</b> Public key $pk$ , message $m \in R_2$ , optional randomness $r$ <b>Output:</b> Ciphertext $c$	
1: $seed_A \leftarrow \mathcal{U}(\{0, 1\}^{256})$		1: $\mathbf{A} \leftarrow \text{gen}(seed_A) \in R_q^{l \times l}$	
2: $\mathbf{A} \leftarrow \text{gen}(seed_A) \in R_q^{l \times l}$		2: <b>if</b> $r$ is not specified <b>then</b>	
3: $r \leftarrow \mathcal{U}(\{0, 1\}^{256})$		3: $r \leftarrow \mathcal{U}(\{0, 1\}^{256})$	
4: $\mathbf{s} \leftarrow \beta_\mu(R_q^{l \times 1}; r)$		4: $\mathbf{s}' \leftarrow \beta_\mu(R_q^{l \times 1}; r)$	
5: $\mathbf{b} \leftarrow ((\mathbf{A}^T \mathbf{s} + \mathbf{h}) \bmod q) \gg (\epsilon_q - \epsilon_p) \in R_p^{l \times 1}$		5: $\mathbf{b}' \leftarrow ((\mathbf{A} \mathbf{s}' + \mathbf{h}) \bmod q) \gg (\epsilon_q - \epsilon_p) \in R_p^{l \times 1}$	
6: <b>return</b> ( $pk := (seed_A, \mathbf{b}), sk := \mathbf{s}$ )		6: $v' \leftarrow \mathbf{b}'^T(\mathbf{s}' \bmod p) \in R_p$	
		7: $c_m \leftarrow (v' + h_1 - 2^{\epsilon_p - 1} m \bmod p) \gg (\epsilon_p - \epsilon_T) \in R_T$	
		8: <b>return</b> $c := (c_m, \mathbf{b}')$	
Algorithm	II.2		
Saber.PKE.Dec( $sk, c$ )			
<b>Input:</b> Secret key $sk$ , ciphertext $c$ <b>Output:</b> Message $m'$			
1: $v \leftarrow \mathbf{b}'^T(\mathbf{s} \bmod p) \in R_p$			
2: $m' \leftarrow ((v - 2^{\epsilon_p - \epsilon_T} c_m + h_2) \bmod p) \gg (\epsilon_p - 1) \in R_2$			
3: <b>return</b> $m'$			

### B. FrodoKEM

FrodoKEM [2] is a lattice-based KEM that relies on the hardness of Learning With Errors. The submission to NIST specifies the parameters as in the table below.

Parameter set	Sec. level	$n$	$q$	$\bar{m} = \bar{n}$	$l_A$	$l_{SE}$
Frodo-640	1	640	$2^{15}$	8	128	128
Frodo-976	3	976	$2^{16}$	8	128	192
Frodo-1344	5	1344	$2^{16}$	8	128	256

FrodoKEM uses functions  $\text{Gen}(s)$ , to generate a pseudorandom matrix  $\mathbf{A} \in \mathbb{Z}_q^{n \times n}$  from a seed  $s$  of length  $l_A$  (using AES or SHAKE), and  $\text{SM}(r, s, t)$  for inversion sampling of a matrix in  $\mathbb{Z}_q^{s \times t}$  using a pseudorandom array of 16-bit integers  $\mathbf{r}$  and a precomputed table  $T_\chi$  for error distribution  $\chi$ . Let SK denote SHAKE. The PKE is specified by Algs. II.4, II.6, II.5.

### C. The AMX coprocessor

AMX is a matrix multiplication coprocessor found in Apple SoCs. It lacks official documentation, so we turn to the reverse engineering efforts of [15]–[17]. We briefly review some concepts and refer to them for more details, as well as the description in [9], on which our algorithmic notation is based.

AMX’s register file is comprised of 80 64-byte registers: 16 input registers, split as 8 X and 8 Y registers, and 64 output Z registers viewed as rows of a matrix, as depicted in Figure 1. Some instructions can address X and Y registers bitwise as 512-byte circular buffers. AMX instructions are encoded within a reserved opcode space of A64; once no longer speculative, the CPU forwards them to the AMX coprocessor.

Data transfer between the CPU is AMX is done solely through memory, using load (`ldx, ldy, ldz`) and store (`stx, sty, stz`) instructions. `extrh` and `extrv` move rows and columns, respectively, of Z to X or Y registers.

---

**Algorithm II.4**  
 FrodoPKE.KeyGen()
 

---

**Input:** None  
**Output:** Key pair  $(pk, sk)$

- 1:  $seed_A \leftarrow \mathcal{U}(\{0, 1\}^{l_A})$
- 2:  $\mathbf{A} \leftarrow \text{Gen}(seed_A)$
- 3:  $seed_{SE} \leftarrow \mathcal{U}(\{0, 1\}^{l_{SE}})$
- 4:  $\mathbf{r} \leftarrow \text{SK}(0 \times 5F || seed_{SE}, 2n\bar{n} \cdot 16)$
- 5:  $\mathbf{S}^T \leftarrow \text{SM}(\mathbf{r}[0 : n\bar{n} - 1], \bar{n}, n)$
- 6:  $\mathbf{E} \leftarrow \text{SM}(\mathbf{r}[n\bar{n} : 2n\bar{n} - 1], n, \bar{n})$
- 7:  $\mathbf{B} = \mathbf{A} + \mathbf{E}$
- 8: **return**  $(pk := (seed_A, \mathbf{B}), sk := \mathbf{S}^T)$

---

**Algorithm II.5**  
 FrodoPKE.Dec(sk, c)
 

---

**Input:** Secret key  $sk$ , ciphertext  $c$   
**Output:** Message  $m'$

- 1:  $\mathbf{M} = \mathbf{C}_2 - \mathbf{C}_1 \mathbf{S}$
- 2: **return**  $m' := \text{Decode}(\mathbf{M})$

---



---

**Algorithm II.6**  
 FrodoPKE.Enc(pk, m)
 

---

**Input:** Public key  $pk$ , message  $m$   
**Output:** Ciphertext  $c$

- 1:  $\mathbf{A} \leftarrow \text{Gen}(seed_A)$
- 2:  $seed_{SE} \leftarrow \mathcal{U}(\{0, 1\}^{l_{SE}})$
- 3:  $\mathbf{r} \leftarrow \text{SK}(0 \times 96 || seed_{SE}, (2\bar{m}\bar{n} + \bar{m}\bar{n}) \cdot 16)$
- 4:  $\mathbf{S}' \leftarrow \text{SM}(\mathbf{r}[0 : \bar{m}\bar{n} - 1], \bar{m}, n)$
- 5:  $\mathbf{E}' \leftarrow \text{SM}(\mathbf{r}[\bar{m}\bar{n} : 2\bar{m}\bar{n} - 1], \bar{m}, n)$
- 6:  $\mathbf{E}'' \leftarrow \text{SM}(\mathbf{r}[2\bar{m}\bar{n} : 2\bar{m}\bar{n} + \bar{m}\bar{n} - 1], \bar{m}, \bar{n})$
- 7:  $\mathbf{B}' = \mathbf{S}' \mathbf{A} + \mathbf{E}'; \mathbf{V} = \mathbf{S}' \mathbf{B} + \mathbf{E}''$
- 8: **return**  $c := (\mathbf{C}_1, \mathbf{C}_2) = (\mathbf{B}', \mathbf{V} + \text{Encode}(m))$

---

	X[0]	⋯	X[n]
Y[0]	Z[0][0] += Y[0]X[0]	⋯	Z[0][n] += Y[0]X[n]
Y[1]	Z[1][0] += Y[1]X[0]	⋯	Z[1][n] += Y[1]X[n]
⋮	⋮	⋱	⋮
Y[n]	Z[n][0] += Y[n]X[0]	⋯	Z[n][n] += Y[n]X[n]

Fig. 1. AMX register file organization.

The vector-mode `mac16` and `vecint` instructions realize vector operations such as addition  $+$  and the Hadamard (pointwise) product  $\circ$ . Outer product of a column by a row vector (the BLAS Level-2 rank-1 update operation `xGER`) is realized by the matrix-mode `mac16` and `matint` instructions. For 16-bit integers, vectors (or matrix rows/columns) are up to 32 elements long; each instruction’s enable modes can mask part of the computation if smaller sizes are needed.

We illustrate the notation with AMX’s primary application, matrix multiplication (in our case,  $32 \times 32$  matrices with 16-bit integer data), in Algorithm II.7. It is also a basic block, with suitable modifications, for our Saber and FrodoKEM AMX implementations. If  $\mathbf{A}^T$  rather than  $\mathbf{A}$  is input to the algorithm, we eschew the transposition by removing lines 1 and 2, and replacing line 4 with a load of the  $i$ -th row of  $\mathbf{A}$ .

---

**Algorithm II.7** MATMULADD( $\mathbf{A}, \mathbf{B}$ ): Compute  $\mathbf{Z}[0 : 2 : 62] \leftarrow \mathbf{Z}[0 : 2 : 62] + \mathbf{A}\mathbf{B}$  using AMX.
 

---

**Input:**  $\mathbf{A}, \mathbf{B} \in \mathbb{Z}_{2^{16}}^{32 \times 32}$  in row-major memory layout.  
**Output:**  $\mathbf{Z}[0 : 2 : 62] + \mathbf{A}\mathbf{B} \in \mathbb{Z}_{2^{16}}^{32 \times 32}$  in even Z registers.

- 1: **for**  $i = 0$  to 31 **do** ▷ Load  $\mathbf{A}$  to odd Z rows
- 2:  $\mathbf{Z}[2i + 1] \leftarrow \text{ldz}(\mathbf{A}[i][0 : 31])$
- 3: **for**  $i = 0$  to 31 **do**
- 4:  $\mathbf{Y}_0 \leftarrow \text{extrv}(\mathbf{Z}[1 : 2 : 63][i])$  ▷  $\mathbf{A}$  transpose step
- 5:  $\mathbf{X}_0 \leftarrow \text{ldx}(\mathbf{B}[i][0 : 31])$
- 6:  $\mathbf{Z}[0 : 2 : 62] \leftarrow \text{mac16}(\mathbf{Z}[0 : 2 : 62] + \mathbf{Y}_0^T \mathbf{X}_0)$

---

We also review the `genlut` instruction, which is instru-

mental to our table-based sampling technique of §IV-D. It has two distinct modes, *generate* and *lookup*. The latter is similar to NEON’s TBL instruction: given an input register with a densely packed array of lane indices (in a format fully described in [17]) and another register containing a table, it performs a table lookup operation; in 16-bit mode, registers are 32 elements wide. The *generate* mode is especially interesting, and unlike any CPU instruction we are familiar with. It takes a table  $T$  and source register  $V$  as input, and generates a packed array of lane indices, in the format used by *lookup* mode, by searching for the minimum index  $i - 1$  such that  $T[i] > V[l]$ , for each lane  $l$  of the source. If  $T$  is sorted in ascending order, `genlut` returns  $i$  such that  $T[i] \leq V[l] < T[i + 1]$ .

We refer to [9] for performance characteristics of AMX; the main improvement reported for M3 is that some vector operations can dual-issue, while M1 is strictly single-issue.

### III. SABER ON AMX

We now discuss AMX-accelerated multiplication in  $\mathbb{Z}_{2^{16}}[X]/(X^n + 1)$ , which is Saber’s main algorithmic task.

#### A. Baseline implementation

An AMX-based algorithm was previously proposed in [9] for multiplication in  $\mathbb{Z}_{2^{16}}[X]/(X^n - 1)$ , which is identical to the multiplication in  $\mathbb{Z}_{2^{16}}[X]/(X^n + 1)$ , except for flipping signs of terms with powers greater than  $n - 1$  in the reduction. We achieve this via `vecint` and `matint` instructions, which generalize vector- and matrix-mode `mac16` (respectively) with accumulation by either adding or subtracting. We refer to [9] for more on the techniques, and only mention the key changes needed to adapt its POLYMODMUL algorithm to Saber:

- Replacing the `mac16` instruction in line 9 of the ACCUMULATEOUTERPRODUCTSREDUCTION subroutine by `matint` using accumulation by subtraction.
- Modifying the `vecint` instructions in lines 8 and 13 of the MERGEFIRSTANDLASTBLOCKS subroutine to perform accumulation by subtraction.

Inner products of polynomials  $\sum_{k=0}^{l-1} a^{(k)}(x)b^{(k)}(x)$  can be computed naively by  $l$  POLYMODMUL calls; for each  $k$  and each 32-coefficient slice of the output polynomial, ACCUMULATEOUTERPRODUCTSREDUCTION sums terms of the form  $\mathbf{b}_j^{(k)T} \mathbf{a}_i^{(k)}$ . We propose accumulating  $\sum_{k=0}^{l-1} \mathbf{b}_j^{(k)T} \mathbf{a}_i^{(k)}$  instead, lazily applying POLYMODMUL’s vector operations (shifts and flattenings) for a factor- $l$  increase in the key AMX performance metric (matrix/vector operation ratio).

Encryption and key generation multiply an  $l \times l$  polynomial matrix  $\mathbf{A}$  by an  $l \times 1$  polynomial vector  $\mathbf{s}$ ; we compute these as inner products between rows of  $\mathbf{A}$  and  $\mathbf{s}$ . This is followed by element-wise addition of a constant polynomial  $\mathbf{h}$  and right-shifting by a constant  $\epsilon_q - \epsilon_p$ . Both tasks can be accelerated by AMX, by loading  $\mathbf{h}$  to the Z registers before accumulating the matrix-vector product and using a specific ALU mode of `vecint` for shifting (see Algorithm III.3 for an example of how these steps are incorporated). In Saber’s specification, all coefficients of  $\mathbf{h}$  are identical, so they can be distributed to all Z registers with a single `mac16` instruction.

## B. TMVP-based implementation

We now present a second method for polynomial multiplication, based on the *Toeplitz matrix-vector product*, which has the option of accumulating to the Z registers the batched multiplication of a single polynomial  $b(x)$  by multiple polynomials  $a^{(l)}(x)$ . This has potential to improve the matrix-vector products in Saber’s encryption and key generation which, using LightSaber ( $l = 2$ ) as an example, can be written as

$$\mathbf{As} = \begin{pmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{pmatrix} \begin{pmatrix} s_0 \\ s_1 \end{pmatrix} = s_0 \begin{pmatrix} A_{00} \\ A_{10} \end{pmatrix} + s_1 \begin{pmatrix} A_{01} \\ A_{11} \end{pmatrix}, \quad (1)$$

so that multiplications are performed in two batches of two with the additions coming free by accumulation in Z registers.

In the TMVP approach, the coefficients for a single product  $c(x) := a(x)b(x)$  are computed as

$$\begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ \vdots \\ c_{n-1} \end{pmatrix} = \begin{pmatrix} b_0 & -b_{n-1} & -b_{n-2} & \cdots & -b_1 \\ b_1 & b_0 & -b_{n-1} & \cdots & -b_2 \\ qb_2 & b_1 & b_0 & \cdots & -b_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{n-1} & b_{n-2} & b_{n-3} & \cdots & b_0 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{pmatrix},$$

which we denote  $\mathbf{c} = \mathbf{M}\mathbf{a}$ . We represent the batched products by promoting  $\mathbf{c}$  and  $\mathbf{a}$  to matrices, with one column per  $a^{(l)}(x)$ , a trick first used in the CUDA implementation from [18]. For the remainder, we fix  $n = 256$  and assume for illustration purposes that only two multiplications are batched (this will be the case in LightSaber). By splitting  $\mathbf{M}$  into  $32 \times 32$  blocks and each of the  $\mathbf{c}^{(l)}, \mathbf{a}^{(l)}$  into  $32 \times 1$  blocks, we get

$$\begin{pmatrix} C_0^{(0)} & C_0^{(1)} \\ C_1^{(0)} & C_1^{(1)} \\ \vdots & \vdots \\ C_7^{(0)} & C_7^{(1)} \end{pmatrix} = \begin{pmatrix} B_0 & -B_7 & \cdots & -B_2 & -B_1 \\ B_1 & B_0 & \cdots & -B_3 & -B_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ B_7 & B_6 & \cdots & B_1 & B_0 \end{pmatrix} \begin{pmatrix} A_0^{(0)} & A_0^{(1)} \\ A_1^{(0)} & A_1^{(1)} \\ \vdots & \vdots \\ A_7^{(0)} & A_7^{(1)} \end{pmatrix},$$

and by exploiting the Toeplitz (more precisely, skew-circulant) property retained by the  $B_i$ , this can be rearranged to

$$\begin{pmatrix} C_0^{(0)} & C_0^{(1)} & C_1^{(0)} & C_1^{(1)} & \cdots & C_7^{(0)} & C_7^{(1)} \end{pmatrix} = \begin{pmatrix} B_0 & A_0^{(0)} & A_0^{(1)} & A_1^{(0)} & A_1^{(1)} & \cdots & A_7^{(0)} & A_7^{(1)} \\ +B_1 & -A_7^{(0)} & -A_7^{(1)} & A_0^{(0)} & A_0^{(1)} & \cdots & A_6^{(0)} & A_6^{(1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ +B_7 & -A_1^{(0)} & -A_1^{(1)} & -A_2^{(0)} & -A_2^{(1)} & \cdots & A_0^{(0)} & A_0^{(1)} \end{pmatrix} \quad (2)$$

which we denote  $\sum_{i=0}^7 B_i \mathcal{A}_i$ , defining the  $32 \times 16$  matrices in parentheses as  $\mathcal{A}_i$ . Note that each  $\mathcal{A}_i$  can be obtained from different 16-element-wide slices of the  $32 \times 30$  matrix

$$\mathbf{A} := \begin{pmatrix} -A_1^{(0)} & -A_1^{(1)} & \cdots & -A_7^{(0)} & -A_7^{(1)} & A_0^{(0)} & A_0^{(1)} & \cdots & A_7^{(0)} & A_7^{(1)} \end{pmatrix}. \quad (3)$$

Algorithm III.1 stores the transpose of this matrix (since it is more efficient to load coefficient slices into rows) to the 30 largest odd-numbered Z registers. Likewise, the  $B_i$  matrices for  $i = 0$  and  $i > 0$  are given respectively by

$$B_0 = \begin{pmatrix} b_0 & -b_{255} & \cdots & -b_{225} \\ b_1 & b_0 & \cdots & -b_{226} \\ \vdots & \vdots & \ddots & \vdots \\ b_{31} & b_{30} & \cdots & b_0 \end{pmatrix}, \quad B_i = \begin{pmatrix} b_{32i} & b_{32i-1} & \cdots & b_{32i-31} \\ b_{32i+1} & b_{32i} & \cdots & b_{32i-30} \\ \vdots & \vdots & \ddots & \vdots \\ b_{32i+31} & b_{32i+30} & \cdots & b_{32i} \end{pmatrix},$$

so every column of every matrix can be obtained from a 32-element-tall slice of the 287-element column vector

$$\begin{pmatrix} -b_{225} & -b_{226} & \cdots & -b_{255} & b_0 & b_1 & \cdots & b_{255} \end{pmatrix}^T.$$

Note that all but the negative terms (used only for  $B_0$ ) fit in the X registers, so we load only  $(b_0 \cdots b_{255})$  to those registers and replace  $b_{224}, \dots, b_{255}$  by their negated version as needed.

Our algorithm works with the transpose of (2), so output coefficients can be stored to memory in the natural row-major layout; thus, we compute  $\sum_{i=0}^7 \mathcal{A}_i^T B_i^T = \sum_{i=0}^7 \sum_{j=0}^{31} \mathcal{A}_i^T[:, j] B_i^T[j, :]$ . Here,  $\mathcal{A}_i^T[:, j] B_i^T[j, :]$  corresponds to the outer product of  $X[32i-j : 32i-j+31]$  and  $Z[33-4i : 2 : 63-4i][j]$  (with the latter obtainable via an `extrv` instruction). The resulting algorithm is presented as Algorithm III.2.

*Remark 1:* It is straightforward to generalize the method in this section to the case of batching  $l$  polynomial multiplications of  $\mathbf{b}(x)$  by  $\mathbf{a}^{(0)}(x), \dots, \mathbf{a}^{(l-1)}(x)$ ; the  $B_i$  remain the same whereas the  $\mathcal{A}_i$  become matrices of dimension  $32 \times 8l$ . The outer products grow to size  $32 \times 8l$  and can still be computed with a single `mac16` instruction as long as  $l \leq 4$  (which covers all Saber parameter sets). Meanwhile, the matrix  $\mathcal{A}^T$  becomes of size  $15l \times 32$ , so it is no longer possible to store it in one half of the Z registers for  $l > 2$ . Instead, one can modify Algorithm III.1 and Algorithm III.2 to spill and reload rows of  $\mathcal{A}^T$  on demand using an external array, introducing some overhead due to AMX loads and stores.

Algorithm III.3 covers the entire computation of  $(\mathbf{As} + \mathbf{h}) \gg (\epsilon_q - \epsilon_p)$  for LightSaber, with straightforward generalizations for Saber ( $l = 3$ ) and FireSaber ( $l = 4$ ) per Remark 1.

---

**Algorithm III.1** PREPAREMATRIXA( $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}$ ): Loads  $\mathcal{A}^T$  from equation (3) to odd Z registers.

---

**Input:**  $\mathbf{a}^{(0)}$  and  $\mathbf{a}^{(1)}$  (arrays of 256 coefficients each)

**Output:** Loads  $-a_{32i:32i+31}^{(l)}$  to  $Z[2(2i+1)+1]$  for  $0 < i \leq 7$ , and  $a_{32i:32i}^{(l)}$  to  $Z[2(2i+1+16)+1]$  for  $0 \leq i \leq 7$

- 1:  $Y_0 \leftarrow \text{ldy}([-1, \dots, -1])$
- 2: **for**  $l = 0$  to 1 **do**
- 3:     **for**  $i = 0$  to 7 **do**
- 4:          $Z[2(2i+l+16)+1] \leftarrow \text{ldz}(\mathbf{a}^{(l)}[32i : 32i+31])$
- 5:          $X_0 \leftarrow \text{ldx}(\mathbf{a}^{(l)}[32i : 32i+31])$
- 6:          $Z[2(2i+l)+1] \leftarrow \text{mac16}(X_0 \circ Y_0)$

---

## IV. FRODOKEM ON NEON AND AMX

In this section, we discuss implementation techniques to speed up FrodoKEM using NEON, and then our AMX implementation which achieves further significant speedups.

### A. NEON optimizations

The main improvements in our NEON implementation come from (i) refining the generation of the matrix  $\mathbf{A}$  in the AES variant and (ii) a careful loop order for matrix multiplication.

The reference implementation initializes the “plaintext” for  $\mathbf{A}$  in a first pass and encrypts it with AES in a second pass. We do it in a single pass, generating the “plaintext” in NEON registers and encrypting with ARMv8 AES instructions.

---

**Algorithm III.2** POLYMODMULTMVP( $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \mathbf{b}$ ): Multiplication in  $\mathbb{Z}_{2^{16}}[X]/(X^{256} + 1)$  of a polynomial  $\mathbf{b}$  by two polynomials  $\mathbf{a}^{(l)}$  using AMX.

---

**Input:**  $\mathbf{b}, \mathbf{a}^{(0)}, \mathbf{a}^{(1)}$  (arrays of 256 coefficients)

**Output:** Accumulates to  $Z[0 : 2 : 30]$  the coefficients for  $c^{(l)}(x) = a^{(l)}(x)b(x)$ , mapping  $c_{32j:32j+31}^{(l)}$  to  $Z[4j + 2l]$ .

- 1:  $X_0, \dots, X_7 \leftarrow \text{ldx}(\mathbf{b}[0 : 31]), \dots, \text{ldx}(\mathbf{b}[224 : 255])$
- 2:  $\text{tmp} \leftarrow \text{stz}(\text{mac16}(X_7 \circ [-1, \dots, -1]))$
- 3: PrepareMatrixA( $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}$ )  $\triangleright$  load  $\mathcal{A}^T$  to odd  $Z$
- 4: **for**  $j = 0$  **to** 31 **do**
- 5:      $Y_0 \leftarrow \text{extrv}(Z[1 : 2 : 63][j])$   $\triangleright$  extract  $\mathcal{A}^T[:, j]$
- 6:     **for**  $i = 0$  **to** 7 **do**
- 7:         **if**  $i == 0$ :  $X_7 \leftarrow \text{tmp}$   $\triangleright$  negate  $[b_{224}, \dots, b_{255}]$
- 8:          $Z[0 : 2 : 30] \leftarrow \text{mac16}(Z[0 : 2 : 30] + Y[16 - 2i : 31 - 2i]^T X[32i - j : 32i - j + 31])$
- 9:         **if**  $i == 0$ :  $X_7 \leftarrow \text{ldx}(\mathbf{b}[224 : 255])$   $\triangleright$  restore

---

**Algorithm III.3** SABERMATVECMULTMVP( $\mathbf{c}, \mathbf{A}, \mathbf{b}, \mathbf{h}, \epsilon$ ): Computes  $(\mathbf{A}\mathbf{b} + \mathbf{h}) \gg \epsilon$  in  $\mathbb{Z}_{2^{16}}[X]/(X^{256} + 1)$ , shifting coefficient-wise, for  $\mathbf{A}$  a  $2 \times 2$  polynomial matrix and  $\mathbf{b}, \mathbf{h}$   $2 \times 1$  polynomial vectors with all coefficients in  $\mathbf{h}$  equal.

**Input:**  $h$  (repeating coefficient of  $\mathbf{h}$ ),  $\mathbf{b}$  ( $2 \times 256$  coefficient array),  $\mathbf{A}$  ( $2 \times 2 \times 256$  coefficient array),  $\epsilon$  (integer).

**Output:**  $\mathbf{c}$  (array of 256 coefficients for the result)

- 1:  $Z[0 : 2 : 62] \leftarrow \text{mac16}(\text{ldx}([h, h, \dots, h])) \triangleright$  copies of  $h$
- 2: **for**  $i = 0, 1$ : POLYMODMULTMVP( $\mathbf{b}[i], \mathbf{A}[0][i], \mathbf{A}[1][i]$ )
- 3: **for**  $i = 0, 1$  and  $j = 0, \dots, 7$  **do**
- 4:      $\mathbf{c}[i][32j : 32j + 31] \leftarrow \text{stz}(\text{vecint}(Z[4j + 2i] \gg \epsilon))$

---

Let  $\mathbf{U}$  and  $\mathbf{V}$  be matrices of size  $m \times n$  and  $n \times p$ , respectively. Then,  $t_{i,j} = \sum_{k=1}^n u_{i,k}v_{k,j}$  is the entry in row  $i$  and column  $j$  of  $\mathbf{T} = \mathbf{UV}$ . Thus, computing  $\mathbf{T}$  requires three nested `for` loops, iterating through the values of  $i, j$  and  $k$ . Although the order of the loops is arbitrary, data access patterns differ. After evaluating all loop orders, we implemented the one maximizing arithmetic intensity for each of FrodoKEM’s matrix multiplication routines. Performance counters show that our choices below yield  $> 90\%$  ALU usage.

Implementation	Operation			
	$\mathbf{AS} + \mathbf{E}$	$\mathbf{S}'\mathbf{A} + \mathbf{E}$	$\mathbf{B}'\mathbf{S}$	$\mathbf{S}'\mathbf{B} + \mathbf{E}$
Reference	$ijk$	$jik$	$ijk$	$ijk$
Optimized	$ijk$	$kij$	$ijk$	$ijk$
NEON	$ijk$	$kji$	$ikj$	$kij$

Also, for a potential doubling in throughput on the M1 and M3, we employ multiply-accumulate instructions instead of separate multiplication and addition instructions as in [13].

### B. Matrix multiplication on AMX

We focus on the computations  $\mathbf{AS} + \mathbf{E}$  and  $\mathbf{S}'\mathbf{A} + \mathbf{E}'$ , where  $\mathbf{A}$  has size  $n \times n$ ,  $\mathbf{S}, \mathbf{E}$  have size  $n \times \bar{n}$  and  $\mathbf{S}', \mathbf{E}'$  have size  $\bar{n} \times n$ . For both multiplications,  $\mathbf{AS}$  and  $\mathbf{S}'\mathbf{A}$ , the main idea is to load  $\mathbf{A}$  and  $\mathbf{S}$  (or  $\mathbf{S}'$ ) into the  $X$  and  $Y$  input registers,

and compute multiply-accumulates to the  $Z$  output registers. By initializing  $Z$  with  $\mathbf{E}$  or  $\mathbf{E}'$ , we obtain addition “for free”.

In AMX, the natural size of matrices for multiplication is up to  $32 \times 32$ . So,  $\mathbf{AS}$  and  $\mathbf{S}'\mathbf{A}$  are computed via block matrix multiplication. We decompose  $\mathbf{A}$  into  $32 \times 32$  blocks,  $\mathbf{S}$  into  $32 \times 8$  blocks and  $\mathbf{S}'$  into  $8 \times 32$  blocks. However, one of the dimensions being  $< 32$  is sub-optimal since the remaining 24 rows/columns are masked out, but `mac16` does not appear to execute any faster. This wasted computational power is reclaimed through batching in §IV-C. For Frodo-976,  $32 \nmid n$ , so part of the computation for blocks at the edges is masked out, thus behaving as if they were padded with zeros.

Recall that AMX multiplies matrices via outer products of vectors. To compute  $\mathbf{AS}$ , we read blocks of  $\mathbf{A}$  by columns and  $\mathbf{S}$  by rows.  $\mathbf{A}$  and  $\mathbf{S}^T$  are generated in row-major order by Algorithm II.4 (lines 2 and 5); thus,  $\mathbf{S}$  is in column-major order, and we must transpose both  $\mathbf{A}$  and  $\mathbf{S}^T$ . For  $\mathbf{S}'\mathbf{A}$ , both matrices are generated in row-major order by Algorithm II.6 (lines 1 and 4). Thus, we transpose  $\mathbf{S}'$  only.

We report on transposition strategies that performed best among all that we tried. For  $\mathbf{AS}$ , we first transpose the full  $\mathbf{S}^T$  directly in  $C$  (which the compiler autovectorizes to NEON), while  $\mathbf{A}$  is transposed online with AMX during multiplication (as in Algorithm II.7), after generating 32 rows with our one-pass strategy of §IV-A; this is shown in Algorithm IV.1. For  $\mathbf{S}'\mathbf{A}$ , AMX transposition of  $\mathbf{A}$  also performs best.

---

**Algorithm IV.1** FRODO-AS-PLUS-E-32ROWS( $C, \bar{\mathbf{A}}, \mathbf{S}^T, \mathbf{E}, r$ ): Computes rows  $r, \dots, r + 31$  of  $C \leftarrow \mathbf{AS} + \mathbf{E}$ ;  $\bar{\mathbf{A}}$  is the submatrix of  $\mathbf{A}$  containing rows  $r, \dots, r + 31$ .

---

**Input:**  $\bar{\mathbf{A}} \in \mathbb{Z}_{2^{16}}^{32 \times n}$ ;  $\mathbf{S}^T, \mathbf{E} \in \mathbb{Z}_{2^{16}}^{n \times \bar{n}}$ ;  $r \in \{0, 32, \dots, n - 32\}$

**Output:**  $C \in \mathbb{Z}_{2^{16}}^{n \times \bar{n}}$  with rows  $r, \dots, r + 31$  updated.

- 1: Load  $\mathbf{E}[r : r + 31]$  to  $Z[0 : 2 : 62]$
- 2: **for**  $j_0 = 0, 32, 64, \dots, n - 32$  **do**
- 3:     Load  $\mathbf{S}^T[j_0][0 : 7] \parallel \dots \parallel \mathbf{S}^T[j_0 + 31][0 : 7]$  to  $X$
- 4:     Load  $\bar{\mathbf{A}}[0 : 31][j_0 : j_0 + 31]$  to  $Z[1 : 2 : 63]$
- 5:     **for**  $j = 0, \dots, 31$  **do**
- 6:          $Y_0 \leftarrow \text{extrv}(Z[1 : 2 : 63][j])$
- 7:          $Z[0 : 2 : 62][0 : 7] \leftarrow \text{mac16}(Z[0 : 2 : 62][0 : 7] + Y_0^T X[8j : 8j + 7])$
- 8: Store  $Z[0 : 2 : 62]$  to  $C[r : r + 31]$

---

### C. Use of batching

AMX’s throughput is underutilized with single KEM operations as above. We overcome this by introducing an alternate API that batches KEM operations with the same  $(sk, pk)$  pair, i.e., batching multiplications with the same  $\mathbf{A}$ . Thus, it applies to encapsulation and decapsulation (which compute  $\mathbf{S}'\mathbf{A}$ ) but not to key generation (which computes  $\mathbf{AS}$ ).

Batch  $\mathbf{S}'\mathbf{A} + \mathbf{E}$  computation reuses the strategy of §IV-B, except we do 4 computations at once. Recall that  $\mathbf{A}$  has size  $n \times n$  while  $\mathbf{S}'$  and  $\mathbf{E}$  have size  $8 \times n$ . We vertically stack four  $\mathbf{S}'$  or  $\mathbf{E}$  matrices to get  $32 \times n$  matrices, thus fully using AMX’s processing power. The ALUs are nearly saturated

in our NEON implementation (see §IV-A), so batching is implemented straightforwardly (a loop over the 4 copies).

Calculation of  $\mathbf{S}'\mathbf{B} + \mathbf{E}''$  and  $\mathbf{B}'\mathbf{S}$ , without batching, yields small  $8 \times 8$  matrices, too small for profitable use of AMX. With batching, one dimension grows to 32, as in  $\mathbf{AS} + \mathbf{E}$  and  $\mathbf{S}'\mathbf{A}$ , justifying the use of AMX. For  $\mathbf{S}'\mathbf{B} + \mathbf{E}''$ , we perform the single needed transpose online during multiplication as in  $\mathbf{S}'\mathbf{A} + \mathbf{E}$ ; for  $\mathbf{B}'\mathbf{S}$ , which needs two transposes (as  $\mathbf{S}$  is stored transposed in *sk*), we did not achieve a speedup.

#### D. Gaussian sampling using the AMX *genlut* instruction

We present a novel technique for inversion sampling, applied to FrodoKEM Gaussian sampling. At its core is the use of AMX’s *genlut* instruction in generate mode (see §II-C) to perform parallel search on 32 16-bit source values. Its use is straightforward for distributions with non-negative support and tables of  $\leq 31$  elements, and full support if the table fits an X or Y register and inputs use two’s complement representation. The former condition is met by all parameter sets, but for the latter, an incompatible representation (sign-magnitude with the sign given by the least-significant bit) is prescribed. Thus, we must condition the inputs, at some performance cost.

In lieu of actual two’s complement representation, we place the sign at the most significant bit by right-rotating each input (lines 4 and 5 of Algorithm IV.2), and adapt the  $\mathbf{T}_\chi$  tables to work with this format. The algorithm specified in [2] uses a table for the non-negative support only, and applies the sign bit to the output. We avoid separate application of the sign bit in AMX by using two shifted copies of the table. These fit in the table register since the largest  $\mathbf{T}_\chi$  (for Frodo-640) has  $j + 1 = 13$  elements. Concretely, if  $\mathbf{T}_\chi = [t_0, t_1, \dots, t_j]$  is the original table, we map it to the *genlut*-specific table

$$\mathbf{T}'_\chi = [0, t_0 + 1, \dots, t_j + 1, t_0 + 2^{15} + 1, \dots, t_j + 2^{15} + 1].$$

Finally, *genlut* in generate mode outputs a densely packed representation (20 bytes representing the results of 32 parallel searches). The remainder of the FrodoKEM code expects the usual 16 bits per element representation. We use *genlut* in lookup mode to map results to the range  $[-j, j]$ , in accordance with our choice of  $\mathbf{T}'_\chi$ . The 32-element mapping is given by

$$\iota = [0, 1, \dots, j, 0, -1, \dots, -j, -j, \dots, -j].$$

We display this procedure as Algorithm IV.2.

## V. EXPERIMENTAL RESULTS

In this section, we describe our experimental setup, report and analyze performance results, and report on experiments on the constant-time execution of AMX’s *genlut* instruction.

#### A. Experimental setup

We benchmark on Apple laptops (a 2020 MacBook Air with M1 SoC and a 2023 MacBook Pro with M3 Max SoC), running macOS 14 and version 15 of Apple’s clang compiler.

As in [9], we explore distinct array allocation strategies. In the usual stack allocation, neighboring variables of a function are very likely allocated in the same memory page, risking

---

**Algorithm IV.2** SAMPLEMATRIX( $\mathbf{s}, \mathbf{T}'_\chi, \iota$ ):  $s_i \leftarrow \mathbf{T}'_\chi[s_i]$  for  $0 \leq i < \bar{n} \cdot n$ .

---

**Input:**  $\mathbf{s} \in \mathbb{Z}_{2^{16}}^{\bar{n} \times n}$  (uniform samples);  $\mathbf{T}'_\chi, \iota \in \mathbb{Z}_{2^{16}}^{32}$  (as above).

**Output:**  $\mathbf{s} \in \mathbb{Z}_{2^{16}}^{\bar{n} \times n}$  (Gaussian samples)

```

1:  $Y_0, Y_1, Y_2 \leftarrow \text{ldy}(\mathbf{T}'_\chi), \text{ldy}(\iota), \text{ldy}([2^{15}, \dots, 2^{15}])$ 
2: for  $i = 0, 32, 64, \dots, \bar{n} \cdot n - 32$  do Process 32 elements
3:    $X_0 \leftarrow \text{ldx}(\mathbf{s}[i : i + 31])$ 
4:    $Z[0] \leftarrow \text{vecint}(X_0 \circ Y_2)$ 
5:    $Z[0] \leftarrow \text{vecint}(Z[0] + X_0 \ggg 1)$ 
6:    $X[0] \leftarrow \text{extrh}(Z[0])$ 
7:    $X_0 \leftarrow \text{genlut}(\text{mode} = \text{gen}, \text{src} = X_0, \text{tbl} = Y_0)$ 
8:    $X_0 \leftarrow \text{genlut}(\text{mode} = \text{lookup}, \text{src} = X_0, \text{tbl} = Y_1)$ 
9:    $\mathbf{s}[i : i + 31] \leftarrow \text{stx}(X_0)$ 

```

---

concurrent CPU and AMX accesses, which cause performance degradation. The POSIX *mmap*() function returns new memory pages for each allocation, sidestepping this issue.

As symmetric operations make up the bulk of execution time in Saber and FrodoKEM, we use fast implementations of SHAKE, AES and NIST’s *randombytes*() function, of ARMv8’s Cryptographic Extensions. We use the  $2 \times$ -batched SHAKE implementation of Becker et al. [10], and modify their unbatched SHAKE code to use ARMv8’s SHA-3 instructions. We implement AES-ECB (for FrodoKEM) and AES-CTR-DRBG (for *randombytes*()), ensuring outputs match with existing implementations.

We use the macOS cycle counting code of [19], and report an average of 1024 executions for each measurement.

#### B. Performance results

We report Saber and FrodoKEM performance data for KEM operations and specific subroutines accelerated by AMX. For memory allocation strategies (stack or *mmap*()), we pick the fastest for each individual measurement and display it through the background color of each table cell: light gray for stack, white for *mmap*(). We compute speedups as ratios between the previous state-of-the-art implementation and our AMX one; for FrodoKEM, we compare our NEON implementation to the previous state-of-the-art, and our AMX implementation to the fastest CPU implementation (usually, our NEON one).

Saber results are shown in Table I. “MVMR” refers to matrix-vector multiplication of polynomials with rounding. NIST levels 1, 3 and 5 map to LightSaber, Saber and FireSaber parameter sets, respectively. FrodoKEM results without batching are shown in Table II, and with  $4 \times$  batching in Table III. We benchmark KEM operations as well as matrix operations  $\mathbf{AS} + \mathbf{E}$  and  $\mathbf{S}'\mathbf{A} + \mathbf{E}$ . The “full” subheading includes the cost of generating matrix  $\mathbf{A}$ , while “mat. mul.” does not.

In lieu of a full discussion, we highlight the main takeaways:

- Most of the cost in both schemes is for symmetric operations, which use the same implementation everywhere. Amdahl’s law bounds gains due to polynomial/matrix multiplication; results should be viewed in that context.
- An analysis of our full dataset shows that only Saber AMX implementations consistently favor *mmap*() mem-

Sec lvl	Work	Operation							
		Key gen.		Encaps.		Decaps.		MVMR	
		M1	M3	M1	M3	M1	M3	M1	M3
1	[10]	19.0	18.9	26.1	26.0	25.6	25.4	4.02	3.96
	III-A	17.5	17.3	23.9	23.6	22.4	22.0	2.47	2.22
	III-B	17.2	17.2	23.5	23.4	22.1	21.9	2.37	2.32
<b>Speedup(<math>\times</math>)</b>		<b>1.10</b>	<b>1.10</b>	<b>1.11</b>	<b>1.11</b>	<b>1.16</b>	<b>1.16</b>	<b>1.69</b>	<b>1.78</b>
3	[10]	31.3	31.2	40.1	40.0	40.4	40.1	7.60	7.52
	III-A	29.2	28.6	36.6	36.0	35.5	34.7	4.46	4.10
	III-B	28.3	28.2	35.1	35.1	34.1	33.8	3.66	3.62
<b>Speedup(<math>\times</math>)</b>		<b>1.11</b>	<b>1.11</b>	<b>1.14</b>	<b>1.14</b>	<b>1.18</b>	<b>1.18</b>	<b>2.08</b>	<b>2.08</b>
5	[10]	48.3	48.2	59.4	59.3	59.8	59.6	12.3	12.2
	III-A	44.8	44.2	54.1	53.5	52.4	51.7	6.95	6.46
	III-B	42.8	42.7	51.6	51.4	49.9	49.6	4.92	4.83
<b>Speedup(<math>\times</math>)</b>		<b>1.13</b>	<b>1.13</b>	<b>1.15</b>	<b>1.15</b>	<b>1.20</b>	<b>1.20</b>	<b>2.50</b>	<b>2.52</b>

TABLE I  
SABER PERFORMANCE IN KILOCYCLES, COMPARING THE  
IMPLEMENTATION OF [10] TO THE ALGORITHMS OF SEC. III-A AND III-B.

ory allocation (geometric mean over KEM operations and security levels: 13% and 4% faster on M1 and M3, respectively); all other implementations show  $< 1\%$  mean difference between stack and `mmap()`. The FrodoKEM AMX implementation is also unaffected, likely due to its use of large matrices, spanning multiple memory pages.

- Both Saber AMX implementations outperform that of [10], despite replacing NEON’s  $\mathcal{O}(n \log n)$  NTT methods by schoolbook ( $\mathcal{O}(n^2)$ ); the implementation of §III-B is faster than that of §III-A. KEM operations are sped up by 10% to 20%, and matrix-vector multiplication of polynomials by 69% to 152%; gains rise with the parameter  $l$  due to better utilization of AMX, per Remark 1. Decapsulation gains are pronounced as it performs reencryption.
- FrodoKEM’s optimized implementation [12] is up to  $\approx 15\times$  faster than the reference one, reinforcing our belief that it is on par or faster than that of Kwon et al. [13]. Our NEON implementation achieved further speedups of up to 29%, 23% and 29% for key generation, encapsulation and decapsulation, respectively. AMX improves upon NEON by up to 16%, 19% and 21% for these operations.
- Matrix operations are further sped up due to the reduced cost of symmetric operations (even more so when **A** matrix generation is removed). AMX achieves up to 124% gains over our already improved NEON implementation.
- Our AMX-specific Gaussian sampling technique of §IV-D achieves gains of up to 417% over NEON.
- Batching amortizes the cost of generating matrix **A** across 4 operations and ensures full AMX utilization. Our NEON implementation improves upon the optimized one by up to 48%, and AMX improves that further by 90%. For matrix multiplication only, AMX gains up to 709%.

### C. Constant-time behavior of AMX’s `genlut` instruction

Cryptographic code must run in time independent of secret data (or *constant time*) to avoid timing side-channel attacks. Gazzoni Filho et al. [9] verify this for many AMX instructions, but not for `genlut`, used by our sampling method of §IV-D. We remedy this by benchmarking AMX and NEON sampling

routines for different inputs: 0,  $2^{16} - 2$ ,  $2^{16} - 1$  or random inputs. The first three map to specific points in the sampling table: respectively, the first, midpoint and last elements.

We report Frodo-640 results on M3; results for other parameter sets and M1 are similar, and are included in the full results dataset in our GitHub repository. Cycle counts for sampling the full  $8 \times 640$  matrix, across all four inputs, vary from 4586 to 4593 (optimized), 4346 to 4349 (NEON) and 839 cycles (AMX). Thus, modulo small variations across benchmark runs, all implementations appear to run in constant time.

## VI. CONCLUSION AND FUTURE WORK

We have implemented the post-quantum cryptosystems Saber and FrodoKEM using the undocumented AMX CPU-coupled matrix multiplication coprocessor, obtaining considerable speedups over CPU-only implementations.

We highlight the difficulties of fully exploiting AMX’s available processing power. Some strides were made over the work of Gazzoni Filho et al. [9], by recasting Saber polynomial multiplication in matrix-multiplication language; still, only the FireSaber parameter set makes full use of AMX. For FrodoKEM, a batched implementation is needed to achieve this goal. Future cryptosystem designs may wish to revisit parameter choices to favor matrix multiplication accelerators.

We note that the performance of many PQC schemes is dictated by the cost of symmetric operations, rather than arithmetic ones such as polynomial/matrix multiplication. To ensure improvements to the latter are duly reflected in protocol performance, more research is needed (from design, implementation and hardware standpoints) into reducing the share of symmetric operations in the execution time of PQC schemes.

An important class of lattice-based cryptosystems are based on NTTs, such as Kyber and Dilithium; we echo the suggestion of [9] to investigate AMX implementations of such schemes.

Table-based sampling is perceived as difficult to implement efficiently in constant time. Our novel technique of §IV-D, using AMX’s `genlut` instruction, brings renewed hope for such methods; although sampling accounts for a small share of FrodoKEM’s running time, other schemes may benefit considerably more. CPU architects would do well to extend instruction set architectures with a similar instruction.

## ACKNOWLEDGMENTS

The first and last authors thank the CRC – Technology Innovation Institute for funding during the development of this work. We thank the anonymous reviewers for their comments, one of which led to the inner product optimization in §III-A.

## REFERENCES

- [1] J.-P. D’Anvers, A. Karmakar, S. Sinha Roy, and F. Vercauteren, “Saber: Module-LWR based key exchange, CPA-secure encryption and CCA-secure KEM,” in *AFRICACRYPT 2018*, A. Joux, A. Nitaj, and T. Rachidi, Eds. Cham: Springer, 2018, pp. 282–305.
- [2] E. Alkim, J. W. Bos, L. Ducas, P. Longa, I. Mironov, M. Naehrig, V. Nikolaenko, C. Peikert, A. Raghunathan, and D. Stebila, “FrodoKEM learning with errors key encapsulation,” Submission to the NIST PQC Standardization Project, 2021, <https://frodokem.org/files/FrodoKEM-specification-20210604.pdf>.

Sec lvl	Work	Operation															
		Key gen.		Encaps.		Decaps.		AS + E				S'A + E				Gaussian sampling	
		M1	M3	M1	M3	M1	M3	full		mat. mul.		full		mat. mul.		M1	M3
						M1	M3	M1	M3	M1	M3	M1	M3	M1	M3	M1	M3
1	[2]	824	778	7274	6628	7272	6605	581	554	304	296	6917	6252	6636	6022	4.62	4.59
	[12]	624	562	731	670	717	641	398	350	202	192	362	347	186	189	4.62	4.59
	Ours (NEON)	504	465	639	563	623	531	274	266	107	118	280	253	111	108	4.34	4.35
	Ours (AMX)	481	415	582	494	541	449	269	220	79.2	77.6	245	191	52.6	52.4	1.32	0.84
<b>Speedup NEON(<math>\times</math>)</b>	<b>1.24</b>	<b>1.21</b>	<b>1.14</b>	<b>1.19</b>	<b>1.15</b>	<b>1.21</b>	<b>1.45</b>	<b>1.32</b>	<b>1.88</b>	<b>1.63</b>	<b>1.30</b>	<b>1.37</b>	<b>1.68</b>	<b>1.75</b>	<b>1.06</b>	<b>1.05</b>	
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.05</b>	<b>1.12</b>	<b>1.10</b>	<b>1.14</b>	<b>1.15</b>	<b>1.18</b>	<b>1.02</b>	<b>1.21</b>	<b>1.35</b>	<b>1.52</b>	<b>1.14</b>	<b>1.32</b>	<b>2.11</b>	<b>2.05</b>	<b>3.29</b>	<b>5.17</b>	
3	[2]	1777	1689	8885	8796	8834	8727	1438	1349	764	749	8380	8317	7689	7698	6.03	5.99
	[12]	1242	1221	1391	1312	1306	1257	902	877	412	418	854	819	384	383	6.03	5.99
	Ours (NEON)	982	944	1148	1068	1087	1005	642	599	247	247	652	586	256	247	5.64	5.65
	Ours (AMX)	935	843	999	928	928	846	579	528	184	181	635	457	125	125	2.01	1.28
<b>Speedup NEON(<math>\times</math>)</b>	<b>1.27</b>	<b>1.29</b>	<b>1.21</b>	<b>1.23</b>	<b>1.20</b>	<b>1.25</b>	<b>1.41</b>	<b>1.46</b>	<b>1.67</b>	<b>1.69</b>	<b>1.31</b>	<b>1.40</b>	<b>1.50</b>	<b>1.55</b>	<b>1.07</b>	<b>1.06</b>	
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.05</b>	<b>1.12</b>	<b>1.15</b>	<b>1.15</b>	<b>1.17</b>	<b>1.19</b>	<b>1.11</b>	<b>1.13</b>	<b>1.34</b>	<b>1.37</b>	<b>1.03</b>	<b>1.28</b>	<b>2.05</b>	<b>1.98</b>	<b>2.80</b>	<b>4.41</b>	
5	[2]	3023	2878	32418	30619	32338	30579	2561	2429	1309	1282	31702	29841	30460	29469	5.50	5.48
	[12]	2119	1932	2265	2148	2146	2161	1667	1475	822	806	1533	1479	774	766	5.50	5.48
	Ours (NEON)	1662	1610	1916	1766	1839	1680	1204	1124	464	464	1255	1112	512	472	5.08	5.09
	Ours (AMX)	1551	1393	1611	1505	1528	1387	1137	965	334	328	1093	852	229	229	2.78	1.76
<b>Speedup NEON(<math>\times</math>)</b>	<b>1.28</b>	<b>1.20</b>	<b>1.18</b>	<b>1.22</b>	<b>1.17</b>	<b>1.29</b>	<b>1.38</b>	<b>1.31</b>	<b>1.77</b>	<b>1.74</b>	<b>1.22</b>	<b>1.33</b>	<b>1.51</b>	<b>1.62</b>	<b>1.08</b>	<b>1.08</b>	
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.07</b>	<b>1.16</b>	<b>1.19</b>	<b>1.17</b>	<b>1.20</b>	<b>1.21</b>	<b>1.06</b>	<b>1.17</b>	<b>1.39</b>	<b>1.41</b>	<b>1.15</b>	<b>1.30</b>	<b>2.24</b>	<b>2.06</b>	<b>1.83</b>	<b>2.88</b>	

TABLE II  
FRODOKEM-AES PERFORMANCE IN KILOCYCLES, WITHOUT BATCHING. AMX SPEEDUPS RELATIVE TO BEST NEON IMPLEMENTATION.

Sec lvl	Work	Operation							
		Encaps. 4 $\times$		Decaps. 4 $\times$		S'A + E 4 $\times$			
		M1	M3	M1	M3	full		mat. mul.	
						M1	M3	M1	M3
1	[12]	1898	1756	1872	1757	974	911	758	745
	Ours (NEON)	1529	1402	1498	1396	638	615	470	448
	Ours (AMX)	1028	906	994	902	248	209	65.1	63.7
	<b>Speedup NEON(<math>\times</math>)</b>	<b>1.24</b>	<b>1.25</b>	<b>1.25</b>	<b>1.26</b>	<b>1.53</b>	<b>1.48</b>	<b>1.61</b>	<b>1.66</b>
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.49</b>	<b>1.55</b>	<b>1.51</b>	<b>1.55</b>	<b>2.57</b>	<b>2.94</b>	<b>7.22</b>	<b>7.04</b>	
3	[12]	3355	3277	3220	3109	2099	1982	1540	1530
	Ours (NEON)	2688	2598	2535	2420	1482	1422	1087	1001
	Ours (AMX)	1622	1555	1466	1380	543	489	147	147
	<b>Speedup NEON(<math>\times</math>)</b>	<b>1.25</b>	<b>1.26</b>	<b>1.27</b>	<b>1.28</b>	<b>1.42</b>	<b>1.39</b>	<b>1.42</b>	<b>1.53</b>
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.66</b>	<b>1.67</b>	<b>1.73</b>	<b>1.75</b>	<b>2.73</b>	<b>2.91</b>	<b>7.37</b>	<b>6.80</b>	
5	[12]	6319	5931	6187	5705	4399	3945	3184	3067
	Ours (NEON)	4400	4251	4187	4001	2795	2691	2062	1937
	Ours (AMX)	2479	2354	2259	2103	999	884	255	252
	<b>Speedup NEON(<math>\times</math>)</b>	<b>1.44</b>	<b>1.40</b>	<b>1.48</b>	<b>1.43</b>	<b>1.57</b>	<b>1.47</b>	<b>1.54</b>	<b>1.58</b>
<b>Speedup AMX(<math>\times</math>)</b>	<b>1.77</b>	<b>1.81</b>	<b>1.85</b>	<b>1.90</b>	<b>2.80</b>	<b>3.05</b>	<b>8.09</b>	<b>7.68</b>	

TABLE III  
FRODOKEM-AES PERFORMANCE IN KILOCYCLES, WITH 4 $\times$  BATCHING. AMX SPEEDUPS RELATIVE TO BEST NEON IMPLEMENTATION.

[3] BSI, "Migration to post quantum cryptography: Recommendations for action by the BSI," 2021. [Online]. Available: [https://www.bsi.bund.de/SharedDocs/Downloads/EN/BSI/Crypto/Migration\\_to\\_Post\\_Quantum\\_Cryptography.pdf?\\_\\_blob=publicationFile&v=2](https://www.bsi.bund.de/SharedDocs/Downloads/EN/BSI/Crypto/Migration_to_Post_Quantum_Cryptography.pdf?__blob=publicationFile&v=2)

[4] Intl. Organization for Standardization, "FrodoKEM: Learning with errors key encapsulation preliminary draft standard," 2023. [Online]. Available: <https://frodokem.org/files/FrodoKEM-ISO-20230314.pdf>

[5] S. Markidis, S. Chien, E. Laure, I. Peng, and J. S. Vetter, "NVIDIA tensor core programmability, performance & precision," in *2018 IEEE Intl. Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. Los Alamitos, CA, USA: IEEE CompSoc, May 2018, pp. 522–531.

[6] Intel Corporation, "Intel<sup>®</sup> architecture instruction set extensions and future features: Programming reference (revision 047)," December 2022, <https://cdrdv2-public.intel.com/671368/architecture-instruction-set-extensions-programming-reference.pdf>.

[7] F. Wilkinson and S. McIntosh-Smith, "An initial evaluation of Arm's Scalable Matrix Extension," in *2022 IEEE/ACM International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS)*, 2022, pp. 135–140.

[8] A. Rodriguez, *Deep Learning Systems: Algorithms, Compilers, and Processors for Large-Scale Production*, ser. Synthesis Lectures on Computer Architecture. Morgan & Claypool Publishers, Oct. 2020.

[9] D. L. Gazzoni Filho, G. Brandão, and J. López, "Fast polynomial multiplication using matrix multiplication accelerators with applications to NTRU on Apple M1/M3 SoCs," *IACR Communications in Cryptology*, vol. 1, no. 1, 2024.

[10] H. Becker, V. Hwang, M. J. Kannwischer, B.-Y. Yang, and S.-Y. Yang, "Neon NTT: Faster Dilithium, Kyber, and Saber on Cortex-A72 and Apple M1," 2021. [Online]. Available: <https://ia.cr/2021/986>

[11] İ. K. Paksoy and M. Cenk, "TMVP-based multiplication for polynomial quotient rings and application to Saber on ARM Cortex-M4," 2020. [Online]. Available: <https://ia.cr/2020/1302>

[12] J. W. Bos, M. Ofner, J. Renes, T. Schneider, and C. van Vredendaal, "The Matrix Reloaded: Multiplication strategies in FrodoKEM," in *Cryptology and Network Security: 20th International Conference, CANS 2021, Vienna, Austria, December 13-15, 2021, Proceedings*. Berlin, Heidelberg: Springer-Verlag, 2021, p. 72–91.

[13] H. Kwon, K. Jang, H. Kim, H. Kim, M. Sim, S. Eum, W.-K. Lee, and H. Seo, "ARMed Frodo," in *Information Security Applications*, H. Kim, Ed. Cham: Springer, 2021, pp. 206–217.

[14] S. Winograd, *Arithmetic Complexity of Computations*. Society for Industrial and Applied Mathematics, 1980, ch. 6, pp. 57–70.

[15] D. Johnson, "IDA (disassembler) and Hex-Rays (decompiler) plugin for Apple AMX," 2022, <https://gist.github.com/dougallj/7a75a3be1ec69ca550e7c36dc75e0d6f>.

[16] M. Handley, "AArch64-Explore: Exploration of Apple CPUs – volume 3: SoC," 2023, <https://github.com/name99-org/AArch64-Explore>.

[17] P. Cawley, "Apple AMX instr. set," 2023, <https://github.com/corsix/amx/>.

[18] L. Wan, F. Zheng, and J. Lin, "TESLAC: Accelerating lattice-based cryptography with AI accelerator," in *Security and Privacy in Communication Networks*, J. Garcia-Alfaro, S. Li, R. Poovendran, H. Debar, and M. Yung, Eds. Cham: Springer, 2021, pp. 249–269.

[19] D. T. Nguyen and K. Gaj, "Fast NEON-based multiplication for lattice-based NIST post-quantum cryptography finalists," in *Post-Quantum Cryptography*, J. H. Cheon and J.-P. Tillich, Eds. Cham: Springer, 2021, pp. 234–254.