

PHOENIX Crypto-Agile Hardware Sharing for ML-KEM and HQC

Antonio Ras, PhD Candidate

Laboratoire Sécurité des Composants (LSCO), CEA-Leti

Encadrants CEA : Mikael Carmona, Antoine Loiseau, Simon Pontié, Emanuele Valea

Directeurs de thèse : Guénaël Renault, Benjamin Smith

Séminaire SecuElec, Rennes | 25 April 2025



Quantum-threat

- □ Shor algorithm can easily solve hard problem
 - Discrete logarithm (Elliptic Curve Cryptography)
 - □ Integer factorization (RSA)



Quantum-threat

- □ Shor algorithm can easily solve hard problem
 - Discrete logarithm (Elliptic Curve Cryptography)
 - □ Integer factorization (RSA)



Post-Quantum Cryptography

- New matematical hard problem
 - □ <u>Lattice-based</u> and <u>Code-based</u> cryptography □ ...
- □ Key Encapsulation Mechanism (KEM)
- Digital Signature (DS)
- □ End of PQC NIST Process (2017 \rightarrow 2025)



Quantum-threat

- □ Shor algorithm can easily solve hard problem
 - Discrete logarithm (Elliptic Curve Cryptography)
 - □ Integer factorization (RSA)



Quantum-safe transition

- □ **Hybridization** : <u>*Post-quantum*</u> + pre-quantum cryptography
- **Crypto-Agility :** Ability to switch between different PQC solutions



Post-Quantum Cryptography

- New matematical hard problem
 - □ <u>Lattice-based</u> and <u>Code-based</u> cryptography □ ...
- Key Encapsulation Mechanism (KEM)
- Digital Signature (DS)
- □ End of PQC NIST Process (2017 \rightarrow 2025)



Quantum-threat

- Shor algorithm can easily solve hard problem
 - Discrete logarithm (Elliptic Curve Cryptography)
 - □ Integer factorization (RSA)



Quantum-safe transition

- □ **Hybridization** : <u>*Post-quantum*</u> + pre-quantum cryptography
- Crypto-Agility : Ability to switch between different PQC solutions



Post-Quantum Cryptography

- New matematical hard problem
 - Lattice-based and <u>Code-based</u> cryptography
 ...
- Key Encapsulation Mechanism (KEM)
- Digital Signature (DS)
- □ End of PQC NIST Process (2017 → 2025)



Crypto-Agility

- Purpose : Maintain security along PQC families
- **Our targeted agility + Hardware :**
 - □ ML-KEM (lattice-based) + HQC (code-based)

	ML-KEM	HQC	Lattice-code agility			
Hardware implementation	~40	3	2			

Lattice-Code Agility State of the Art

[ZZO+24] Zhu et al, "16.2 A 28nm 69.4kOPS 4.4µJ/Op Versatile Post-Quantum Crypto-Processor Across Multiple Mathematical Problems", 2024 [ZZL+22] Zhu et al, "A 28nm 48KOPS 3.4µJ/Op Agile Crypto-Processor for Post-Quantum Cryptography on Multi-Mathematical Problems", 2022

Agenda

- 1. Crypto-Agility on NIST standards
- 2. Number-Theoretic Transform (NTT)
- **3.** HQC using alternative multiplication
- 4. Proposed sharing strategy
- **5.** Integration results
- 6. Conclusion

Crypto-Agility on NIST standards

A case study for ML-KEM and HQC

Targeted Crypto-agility





Targeted Crypto-agility



Targeted Crypto-agility



Targeted Crypto-agility



How to implement Crypto-Agility efficiently ?

Sharing strategy

Software profiling



Sharing strategy

Purpose : Identify common / similar operations

Possible sharing strategies :

Hash function : *straightforward* approach, <u>FIPS 202</u>
 Polynomial multiplication : *approach*

Sharing strategy

Software profiling



Sharing strategy

Purpose : Identify common / similar operations

Possible sharing strategies :

Hash function : *straightforward* approach, <u>FIPS 202</u>
 Polynomial multiplication : *challenge* approach

Which multiplication strategy to mutualize

ML-KEM		HQC				
Number-Theoretic Transform	&	Sparse - Dense : Not constant time 2-way Karatsuba : New bottleneck Frobenius Additive FFT : Our proposal				

Sharing strategy

Software profiling



Sharing strategy

Purpose : Identify common / similar operations

Possible sharing strategies :

Hash function : *straightforward* approach, <u>FIPS 202</u>
 Polynomial multiplication : <u>Challener</u> approach

Which multiplication strategy to <u>mutualize</u>

ML-KEM		HQC
Number-Theoretic Transform	&	Sparse - Dense : Not constant time 2-way Karatsuba : New bottleneck Frobenius Additive FFT : Our proposal

What about hardware Crypto-Agility ?

System on Chip and hardware agility

System-on-Chip



System on Chip and hardware agility

System-on-Chip

□ Application-specific accelerators

- Tightly-coupled : No flexibility, low data exchange latency
- Loosely-coupled : High flexibility, high data exchange latency



System on Chip and hardware agility

System-on-Chip

□ Application-specific accelerators

- Tightly-coupled : No flexibility, low data exchange latency
- Loosely-coupled : High flexibility, high data exchange latency



Hardware resource-sharing



Number-Theoretic Transform

The polynomial multiplication strategy for ML-KEM

FFT-based polynomial multiplication

How to perform multiplication

- □ Evaluating input polynomial in a different domain using an evaluation points set → FFT step
- □ Perform multiplication between evaluated points → **PWM step**
- □ Interpolate back to get final result \rightarrow Inverse FFT (IFFT) step

IMPORTANT : FFT and IFFT are performed using different processing elements, and



> **Number-Theoretic-Transform** : FFT-like for polynomials in $R_q = \frac{\mathbb{Z}_q}{X^{n+1}}$, where q = 3329 and n = 256.



Number-Theoretic-Transform : FFT-like for polynomials in $R_q = \frac{\mathbb{Z}_q}{X^{n+1}}$, where q = 3329 and n = 256.



 \square \mathbb{Z}_q where q = 3329 (12-bits coefficient)



Modular integer multiplication

$$h = (f * g) \mod q$$

Barrett or Montgomery reduction

> **Incomplete NTT** : Factorize modulo $X^n + 1$ into $\frac{n}{2}$ polynomials of *degree* 1.



> **Incomplete NTT** : Factorize modulo $X^n + 1$ into $\frac{n}{2}$ polynomials of *degree* 1.



Hardware State-of-the-Art

- □ Example of fully configurable hardware design
- □ Include CT, GS and PWM butterflies
- □ Same polynomial ring → easy to design





B Frobenius Additive FFT

The alternative polynomial multiplication strategy for HQC

Frobenius Additive FFT (FAFFT)

FFT-like Binary power-of-two polynomial multiplication

- □ New polynomial basis using **Cantor basis**
- □ Encode and Decode : Frobenius map reduce evaluation points



Frobenius Additive FFT (FAFFT)

FFT-like Binary power-of-two polynomial multiplication

- □ New polynomial basis using **Cantor basis**
- □ Encode and Decode : Frobenius map reduce evaluation points



HQC using **FAFFT** in software



FAFFT profiling

- □ Implementation in C, without ASM optimization
- □ Improved version of FAFFT in BIKE [CCK21]
- Adapt for HQC polynomial size
 - □ HQC-128 : 32 768 bits
 - □ HQC-192/256 : 65 536 bits

(%) NOTORIANO 9 1,7 2X BC 2X ENCODE 47,9 23,9 11,9 11,9 1,6 2X FFT PWM IFFT DECODE IBC

FAFFT PROFILING CORTEX-A9 @125 MHZ

HQC using **FAFFT** in software



FAFFT profiling

- □ Implementation in C, without ASM optimization
- □ Improved version of FAFFT in BIKE [CCK21]
- Adapt for HQC polynomial size
 - □ HQC-128 : 32 768 bits
 - □ HQC-192/256 : 65 536 bits

HQC using FAFFT profiling

First FAFFT integration in HQC

- □ Speed up x3.7 / x5 / x8 for respective HQC-X
- **Baseline :** 2-way Karatsuba
- Our proposal : FAFFT



FAFFT PROFILING CORTEX-A9 @125 MHZ

Benchmark HQC-X (Mcycles)



[CCK21] Chen et al, "Optimizing BIKE for the Intel Haswell and ARM Cortex-M4", 2021

Proposed sharing strategy

The solution to manage different butterfly structures with different mathematical foundations

Super-Butterfly

> <u>Challenges</u>

- □ How to combine different butterfly structures ?
 - Different coefficient size

□ How to combine different mathematical foundation ?

- Integer vs Carryless
- Different reduction method

	NTT in ML-KEM	FAFFT in HQC
Structure	Z ₃₃₂₉	$\mathcal{R}_{2^{32}}$
Coefficient size	12-bits	32-bits
Addition	Modulo a - 2220	
Subtraction	100000 q = 3329	Carryless (AOR)
Multiplication	Integer	Carryless
L Reduction	Barrett	Shift-and-Add



Super-Butterfly

Challenges

- □ How to combine different butterfly structures ?
 - Different coefficient size

□ How to combine different mathematical foundation ?

- Integer vs Carryless
- Different reduction method

	NTT in ML-KEM	FAFFT in HQC
Structure	Z ₃₃₂₉	$\mathcal{R}_{2^{32}}$
Coefficient size	12-bits	32-bits
Addition	Modulo a – 2220	
Subtraction	100000 q = 3329	Carryless (AOR)
Multiplication	Integer	Carryless
L Reduction	Barrett	Shift-and-Add







Agility on mathematical foundation

> <u>Agile Modular Arithmetics</u>

Agile Modular Multiplier

Agility on mathematical foundation

> <u>Agile Modular Arithmetics</u>

- \Box Carryless addition \rightarrow XOR operation
- □ Modular addtion and subtraction
 - Two's complement to have only additions

$$c = (a+b) \mod q = \begin{cases} c_t & \text{if } c_t < q \,, \\ c_t + (\overline{q}+1) & \text{otherwise} \end{cases} \qquad \text{where } c_t := a+b \,; \\ d = (a-b) \mod q = \begin{cases} d_t & \text{if } d_t < q \,, \\ d_t + q & \text{otherwise} \end{cases} \qquad \text{where } d_t := a + (\overline{b}+1) \end{cases}$$



Agility on mathematical foundation

Agile Modular Arithmetics

- \Box Carryless addition \rightarrow XOR operation
- Modular addition and subtraction
 - Two's complement to have only additions

$$c = (a+b) \mod q = \begin{cases} c_t & \text{if } c_t < q \,, \\ c_t + (\overline{q}+1) & \text{otherwise} \end{cases} \qquad \text{where } c_t := a+b \,; \\ d = (a-b) \mod q = \begin{cases} d_t & \text{if } d_t < q \,, \\ d_t + q & \text{otherwise} \end{cases} \qquad \text{where } d_t := a + (\overline{b}+1)$$



Agile Modular Multiplier

Hybrid 2-way Karatsuba

□ Schoolbook-based M_x multiplier

$$\begin{aligned} a &= (a_1, a_0) \\ b &= (b_1, b_0) \end{aligned} \qquad \begin{array}{l} z_0 &= \begin{bmatrix} a_0 \times_{M_0} b_0 \\ z_1 &= \begin{bmatrix} a_1 \times_{M_1} b_1 \end{bmatrix} & \text{NTT butterfly multiplication} \\ z_2 &= \begin{bmatrix} (a_1 \oplus a_0) \times_{M_2} (b_1 \oplus b_0) \end{bmatrix} \oplus z_0 \oplus z_1 \end{aligned}$$

 M_0, M_1 : agile multiplier (integer + carryless) M_2 : carryless multiplier



PHOENIX

Features

- Loosely-coupled accelerator
- □ Process Element : Super-Butterfly valorisation
- Coefficient Generation : Generate coefficient pairs to be processed
- Polynomial Memory : Conflict-free memory
 Constant Memory : Different constant values
- □ Control Unit : Orchestrate everything



PHOENIX

Features

- Loosely-coupled accelerator
- Process Element : Super-Butterfly valorisation
- Coefficient Generation : Generate coefficient pairs to be processed
- Polynomial Memory : Conflict-free memory
 Constant Memory : Different constant values
- □ Control Unit : Orchestrate everything

Further optimization for ML-KEM

Vector-like polynomial multiplicationAddition of polynomial error



- → Increase performance
- \rightarrow reducing communication overhead







Crypto-Agility using PHOENIX for all NIST security levels

Resource utilization on FPGA

System-on-Chip integration

- **Zybo-Z7** @125 MHz, Artix-7 FPGA
- □ AXI4-Lite system bus
- □ All NIST security levels



[YMOS21] Bisheh-Niasar et al, "High-Speed NTT-based Polynomial Multiplication Accelerator for CRYSTALS-Kyber Post-Quantum Cryptography", 2021 [IUH22] Itabashi et al, "Efficient modular poly-nomial multiplier for NTT accelerator of Crystals-Kyber", 2022

[LTHW22] Li et al, "Reconfigurable and high-efficiency polynomial multiplication accelerator for Crystals-Kyber.", 2022

[BNAMK21] Bisheh-Niasa et al, "High-speed NTT-based polynomial multiplication accelerator for post-quantum cryptography", 2021

Resource utilization on FPGA

System-on-Chip integration

- □ Zybo-Z7 @125 MHz, Artix-7 FPGA
- □ AXI4-Lite system bus
- □ All NIST security levels

Resource comparison

□ Absence of FAFFT HW implem

Agility bring resource overhead

NTT cycles aligned to state-of-the-art

High-Speed vs Lightweight design

Observation







Performance using PHOENIX

> <u>ML-KEM</u>

- **Speed-up results** $\left(\frac{SW}{SW+HW}\right)$ due to different platforms
- [FSS20] uses tightly-coupled strategy for NTT
- □ Comparing [WZZ+24]
 - Loosely-coupled NTT
 - Assembly optimization for FIPS202
- □ Promising results → system bus overhead

Comparison for ML-KEM 768 (speed-up)



Performance using PHOENIX

> <u>ML-KEM</u>

- **Speed-up results** $\left(\frac{SW}{SW+HW}\right)$ due to different platforms
- □ [FSS20] uses tightly-coupled strategy for NTT
- □ Comparing [WZZ+24]
 - ✤ Loosely-coupled NTT
 - Assembly optimization for FIPS202

□ Promising results → system bus overhead

> <u>HQC</u>

PHOENIX accelerate x2.3 / x2.4 / x2.0 baseline KEM FAFFT-based

Further optimization for HQC (Optimized HQC)

- Same approach conducted in ML-KEM for NTT
- Save some FFT but doubling public key size
- Case scenario with static key
- □ Speed up x1.15 for Encaps / Decaps HQC levels

Comparison for ML-KEM 768 (speed-up)





[FSS20] Fritzmann et al, "RISQ-V: Tightly Coupled RISC-V Accelerators for Post-Quantum Cryptography", 2020 **[WZZ+24]** Wang et al, "Optimized Hardware-Software Co-Design for Kyber and Dilithium on RISC-V SoC FPGA", 2024

Conclusion

Crypto-Agility

- □ The key to maintain PQC-based security
- □ Need of sharing strategy for efficient implementation

Polynomial Multiplication

- □ Alternative Frobenius AFFT improve current HQC bottleneck
- Polynomial multiplication FFT-based for ML-KEM and HQC
- **Challenge 1** : Combine different butterfly structures
- **Challenge 2** : Combine different mathematical foundations
- PHOENIX valorizes Super-Butterfly design as challenges solution

This work proves the existence of efficient sharing strategy on the new NIST KEM standards

Ongoing Improvements

- □ PHOENIX < NTT + FAFFT (?)
- □ ML-KEM result suffers due to system bus overhead



Resource utilization on FPGA



Resource comparison with Hardware NTT

						Т		
	LUT	\mathbf{FF}	BRAM	\mathbf{DSP}	SEC	Cycles M	MHz	μs
This work	5076	4083	8.5	0	3479	236	125	1.88
[YMÖS21]	2543	792	9	4	2935	232	182	1.27
[IUH22]	904	811	2.5	4	1227	268	216	1.24
[LTHW22]	1170	1164	2	4	1638	235	303	0.78
[BNAMK21]	801	717	2	4	1090	324	222	1.46

$$SEC = \frac{LUT}{4} + \frac{FF}{8} + BRAM \times 200 + DSP \times 100$$

FAFFT Background

Another possible way, found in the state-of-the-art, to perform polynomial multiplication in HQC is the AFFT, which is an FFT-based variant adapted to be used in the context of binary field.

Before explaining this technique we have to introduce few mathematical concepts:

- Cantor basis : They are a set of (β_i) satisfies β₀ = 1, β_i² + β_i = β_{i-1} for i > 0. In pratice, it is a new polynomial basis for constructing a finite field as well as an FFT over the field
- Vanishing polynomials : Given a basis (β_i)^{m-1}_{i=0} in the base field, its sequence of subspaces is
 V_i := span{β₀, β₀, ..., β_{i-1}}. Its subspace vanishing polynomials (s_i) are s_i(x) := Π_{a∈Vi}(x - a)
 Properties :
 - (linearity) $s_i(x)$ contains only monomials in the form x^{2^k}
 - (minimal two terms) $s_i(x) = x^{2^i} + x$ iff i is a power of 2
 - (recursivity) $s_i(x) = s_{i-1}^2(x) + s_{i-1}(x)$

Resources utilization of PHOENIX

Module	LUT	\mathbf{FF}	BRAM	DSP
Control Unit	826	131	0	0
Datapath	4250	3952	8.5	0
Coeff memory address	238	128	0	0
Write address	- 30	612	0	0
Polynomial Memories	0	0	8	0
Constant Memories	37	98	0.5	0
Processing element	3319	3114	0	0
1× SuperButterfly	1432	1380	0	0
COMP1	160	0	0	0
COMP2	176	0	0	0
COMP3	948	548	0	0
COMP4	148	0	0	0
$\mathbf{Total} = \mathtt{Control} \ \mathtt{Unit} + \mathtt{Datapath}$	5076	4083	8.5	0

Performance using PHOENIX

		Device	KeyGen	Encaps	Decaps
	[FSS20]	PULPino	939,932 ×1.21	$1,223,887 \times 1.26$	1,051,003 ×1.45
ML-KEM 512	This Work	Zybo-Z7-20	564,189 ×1.14	570,438 ×1.30	639,403 ×1.38
	$[WZZ^+24]$	PolarFire	$326,983 \times 1.91$	415,496 ×2.01	$394,661 \times 2.42$
	[FSS20]	PULPino	$1,768,400 \times 1.19$	$2,138,810 \times 1.23$	$1,889,930 \times 1.36$
ML-KEM 768	This Work	Zybo-Z7-20	916,828 ×1.15	954,992 ×1.25	1,048,867 ×1.33
	$[WZZ^+24]$	PolarFire	536,213 $\times 1.92$	$671,\!082 \times \! 1.98$	639,024 ×2.32
	[FSS20]	PULPino	$2,856,302 \times 1.18$	$3,312,957 \times 1.21$	2,989,896 ×1.32
ML-KEM 1024	This Work	Zybo-Z7-20	1,450,164 ×1.14	1,492,047 ×1.22	1,614,779 ×1.28
	$[WZZ^+24]$	PolarFire	844,008 ×1.89	$1{,}015{,}251\ \times 1.91$	972,598 ×2.18

		Standard HQC					Optimized HQC				
		Ref.	FAFFT		PHOENIX		FAFFT		PHOENIX		
	KeyGen	28.9	7.3	$3.9 \times$	2.8	$10.3 \times$	-	-	-	-	
HQC-128	Encaps	58.7	15.5	$3.8 \times$	6.5	$9.0 \times$	9.6	$6.1 \times$	5.1	$11.5 \times$	
	Decaps	91.9	27.0	$3.4 \times$	13.6	$6.8 \times$	20.6	$4.5 \times$	13.1	7.0 imes	
	KeyGen	85.1	15.8	$5.4 \times$	6.0	$14.2 \times$	-	-	-	-	
HQC-192	Encaps	171.9	33.3	$5.2 \times$	13.7	$12.6 \times$	20.7	$8.3 \times$	10.9	$15.8 \times$	
	Decaps	262.4	54.4	$4.8 \times$	25.2	$10.4 \times$	39.9	6.6 imes	23.7	$11.1 \times$	
HQC-256	KeyGen	155.8	17.7	$8.8 \times$	7.9	$19.7 \times$	-	-	-	-	
	Encaps	314.6	38.4	$8.2 \times$	18.8	$16.7 \times$	25.5	$12.3 \times$	15.6	$20.2 \times$	
	Decaps	483.4	69.1	$7.0 \times$	39.8	$12.1 \times$	52.4	$9.2 \times$	36.1	$13.4 \times$	



HQC : Standard vs Optimized



