

**4th NIST PQC Standardization Conference**

# **Benchmarking and Analysing NIST PQC Lattice-Based Signature Scheme Standards on ARM Cortex M7**

**James Howe**

Senior Research Scientist



# CONTENTS

- 01 Introduction and Motivation**
- 02 Benchmark and Profile Results**
- 03 Constant-Time Issues**

A large, dark sphere with a white wireframe grid is positioned on the left side of the slide. The background is a dark gray with scattered white dots, resembling a starry sky.

# 01

# INTRODUCTION AND MOTIVATION

# What are the PQC standards we have?

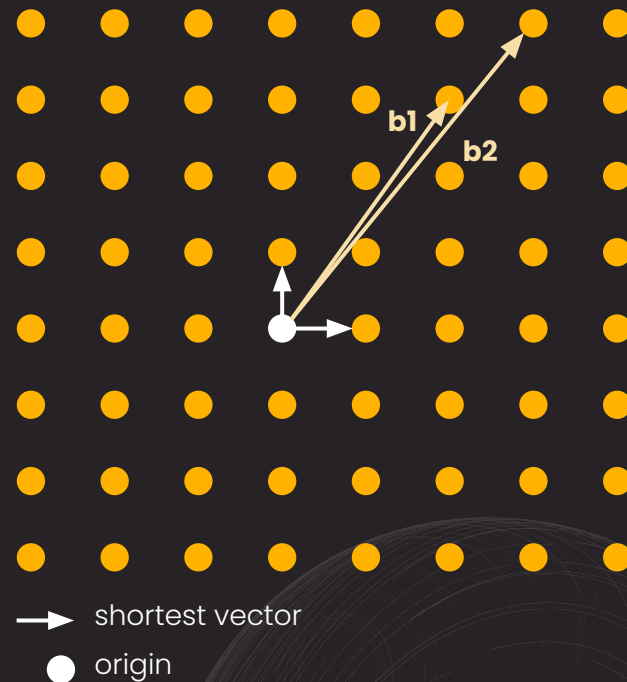
**CRYSTALS-Kyber is the only KEM and CRYSTALS-Dilithium is the primary signature.**

“

“The security of **Kyber** has been thoroughly analyzed [...] based on a strong framework of results in lattice-based cryptography. Kyber has excellent performance overall in software, hardware and many hybrid settings.”

“Dilithium is a signature scheme with high efficiency, relatively simple implementation, a strong theoretical security basis, and an encouraging cryptanalytic history.”

”





# What are the PQC standards we have?

## We also have two other PQ signatures:



**Falcon**, also from lattices, different performance profile.



More complex implementation, emulates or uses FPU.



Offers significantly smaller signature sizes and fast verification.



**Falcon** was chosen for standardization because NIST has confidence in its security (under the assumption that it is correctly implemented) and because its small bandwidth may be necessary in certain applications.

# The Premise

”

“NIST understands that some applications will not work as they are currently designed if the signature and the data being signed cannot fit in a single internet packet.”



“For this reason, NIST decided to standardize FALCON as well. Given FALCON’s overall better performance when signature generation does not need to be performed on constrained devices, many applications may prefer to use FALCON over Dilithium, even in cases in which Dilithium’s signature size would not be a barrier to implementation.”

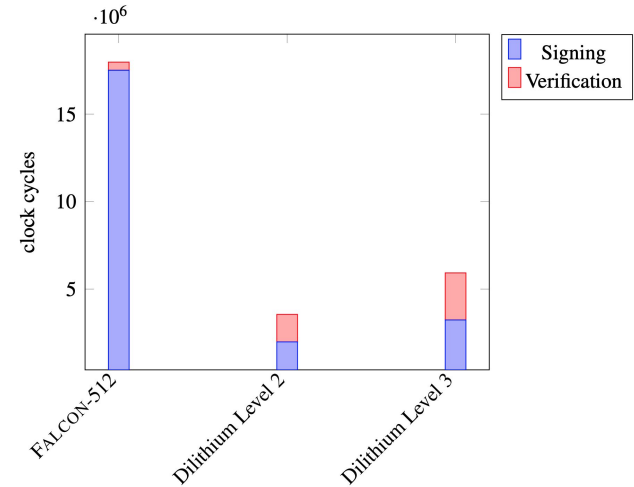


Figure 7. Signature Benchmarks on ARM Cortex-M4 processor

# Current State on ARM Cortex M4

**Without double precision,  
Falcon emulates floats.**

Thus we get **performance profiles** like this on **Cortex M4**.

But can we get this closer using similar device with **full FPU**?

We wanted to **challenge this belief** that Falcon signing is much slower than Dilithium.

**Important decision** in, e.g., RISC-V CPU and SoC implementations.

**Also, does FPU open questions on constant time?**

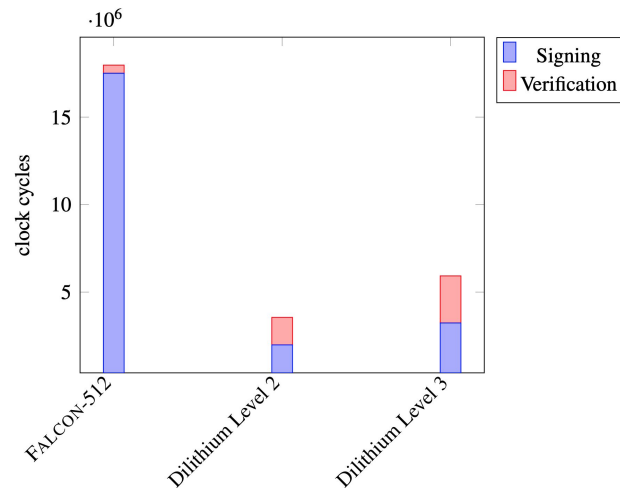


Figure 7. Signature Benchmarks on ARM Cortex-M4 processor

# What's the big deal?

Constant-time and Correctness

01



Emulated **floating-point implementation** can be done

02



Only using integer operations with **uint32\_t** and **uint64\_t** types

03



This is **constant-time**, provided that the underlying platform offers constant-time opcodes for:

- Multiplication of two 32-bit unsigned integers into a 64-bit result.
- Left-shift or right-shift of a 32-bit unsigned integer by a potentially secret shift count in the 0...31 range.

# Why the ARM Cortex M7?

NIST selected Cortex M4 as benchmark MCU;  
and the Cortex M7 is a very similar core

Both have ARMv7-M architecture

Cortex M7 has all ISA features available in the Cortex M4

M7 has 6-stage pipeline (vs 3) and better memory features  
and branch predicting

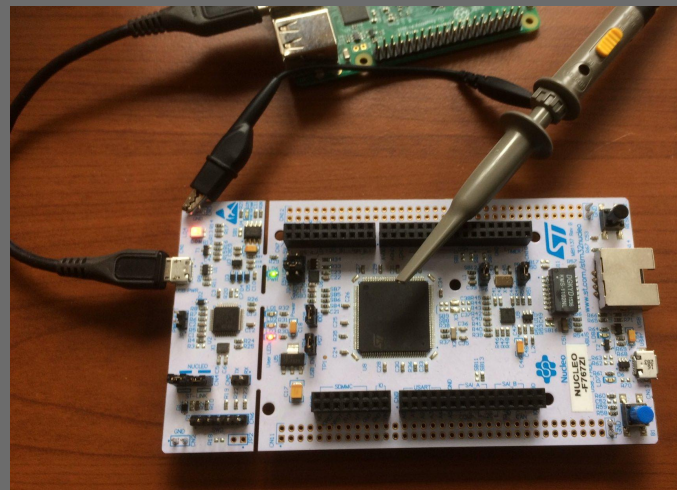


**M7 has 64-bit FPU, M4 has 32-bit**



Falcon requires 53-bit floating-point precision

Using floating-points is rare in cryptography → side channels?



Integer																Floating Point										Cortex-M4 Cortex-M7		Cortex-M4 FPU		Cortex-M7 FPU	
VBBS		VABD		VCMP		VCMPE		VCV		VCVTR		VCY												VFMA		VFP		VFP			
PKHTB		PKHTB		QADD		QADD16		QADD8		QASX		QDADD		QDQUB		SEL		SPLATB		SPLATB		SPLATB		SPLATB							
QASX		QSUB		QSUB16		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8		QSUB8							
SHADD16		SHADD16		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX							
SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX		SHASX							
SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16							
SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16		SHLD16							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC		ADC							
ADC		ADC		ADC																											

A large, dark wireframe sphere is positioned on the left side of the slide, partially obscured by the text. The sphere is composed of many thin, intersecting lines that create a grid-like pattern. The background is a dark, textured surface with subtle light effects and small white specks, resembling a starry sky or a digital environment.

# 02

## **BENCHMARKING AND PROFILING**



# Benchmarking Premise



We benchmarked both Dilithium and Falcon on ARM Cortex M7.



Both used open-source implementations, i.e., pqm4.



Benchmarks took averages over 1000 runs.



All results henceforth are clock cycles, for timings see paper.



We mainly use STM32F767ZI NUCLEO-144 development board.



Using recent GNU ARM embedded toolchain: GCC version 10.2.1 20201103

**using -O2 -mcpu=cortex-m7 -march=-march=armv7e-m+fpv5+fp.dp**

# Dilithium Benchmarking (M4 vs M7)

Overall, the performance of Dilithium wasn't interesting.

Improvements range between 1.09-1.19x

Essentially accounts for the slightly better MCU: Cortex M7 vs the Cortex M4.

Table 1: Benchmarking results of Dilithium on the ARM Cortex M7 using the STM32F767ZI NUCLEO-144 development board. Results in KCycles.

Parameter Set	Operation	Min	Avg	Max	SDev/ SErr	Avg (ms)
Dilithium-2 M7 vs M4	Key Gen	1,390	1,437	1,479	81/3	6.7
	Key Gen	1.13x	<b>1.10x</b>	1.06x	-/-	<b>1.40x</b>
Dilithium-2 M7 vs M4	Sign	1,835	3,658	16,440	604/17	16.9
	Sign	1.19x	<b>1.09x</b>	0.64x	-/-	<b>1.40x</b>
Dilithium-2 M7 vs M4	Verify	1,428	1,429	1,432	27.8/0.9	6.6
	Verify	1.12x	<b>1.12x</b>	1.12x	-/-	<b>1.42x</b>
Dilithium-3 M7 vs M4	Key Gen	2,563	2,566	2,569	37.6/1.2	11.9
	Key Gen	1.12x	<b>1.13x</b>	1.12x	-/-	<b>1.44x</b>
Dilithium-3 M7 vs M4	Sign	2,981	6,009	26,208	65/9	20.7
	Sign	1.12x	<b>1.19x</b>	0.78x	-/-	<b>2.06x</b>
Dilithium-3 M7 vs M4	Verify	2,452	2,453	2,456	26.5/0.8	11.4
	Verify	1.12x	<b>1.12x</b>	1.11x	-/-	<b>1.43x</b>
Dilithium-5	KeyGen	4,312	4,368	4,436	54.4/1.7	20.2
Dilithium-5	Sign	5,020	8,157	35,653	99k/3k	37.8
Dilithium-5	Verify	4,282	4,287	4,292	46.5/1.5	19.8

# Benchmarking Results (FPU vs EMU on M7)



Falcon sees a drastic speedup, expectedly



Improvements range between >6-8x overall



Key generation is least impacted, >1.5x speedup overall.



Signing times show most improvements:

- Sign dynamic >6x speedup, close to Dilithium performance.
- Sign tree >4.5x speedup, comfortably faster than Dilithium

**Verify not impacted, doesn't require floats.**

Table 2: Benchmarking results of Falcon on the ARM Cortex M7 using the STM32F767ZI NUCLEO-144 development board. Results in Kcycles.

Parameter Set	Operation	Min	Avg	Max	SDev/SErr	Avg (ms)
Falcon-512-FPU	Key Gen	44,196	77,475	256,115	226k/7k	358.7
Falcon-512-EMU	Key Gen	76,809	128,960	407,855	303k/9k	597.0
FPU vs EMU	Key Gen	1.74x	<b>1.66x</b>	1.59x	-/-	<b>1.66x</b>
Falcon-1024-FPU	Key Gen	127,602	193,707	807,321	921k/29k	896.8
Falcon-1024-EMU	Key Gen	202,216	342,533	1,669,083	2.4m/76k	1585.8
FPU vs EMU	Key Gen	1.58x	<b>1.76x</b>	2.07x	-/-	<b>1.77x</b>
Falcon-512-FPU	Sign Dyn	4,705	4,778	4,863	149/4	22.1
Falcon-512-EMU	Sign Dyn	29,278	29,447	29,640	188/6	136.3
FPU vs EMU	Sign Dyn	6.22x	<b>6.16x</b>	6.10x	-/-	<b>6.17x</b>
Falcon-1024-FPU	Sign Dyn	10,144	10,243	10,361	1408/44	47.4
Falcon-1024-EMU	Sign Dyn	64,445	64,681	64,957	3k/101	299.5
FPU vs EMU	Sign Dyn	6.35x	<b>6.31x</b>	6.27x	-/-	<b>6.32x</b>
Falcon-512-FPU	Sign Tree	2,756	2,836	2,927	6/.2	13.1
Falcon-512-EMU	Sign Tree	13,122	13,298	13,506	126/4	61.6
FPU vs EMU	Sign Tree	4.76x	<b>4.69x</b>	4.61x	-/-	<b>4.70x</b>
Falcon-1024-FPU	Sign Tree	5,707	5,812	5,919	1422/45	26.9
Falcon-1024-EMU	Sign Tree	28,384	28,621	28,877	3k/115	132.5
FPU vs EMU	Sign Tree	4.97x	<b>4.92x</b>	4.88x	-/-	<b>4.93x</b>
Falcon-512-FPU	Exp SK	1,406	1,407	1,410	8.6/0.3	6.5
Falcon-512-EMU	Exp SK	11,779	11,781	11,788	7/0.2	54.5
FPU vs EMU	Exp SK	8.38x	<b>8.37x</b>	8.36x	-/-	<b>8.38x</b>
Falcon-1024-FPU	Exp SK	3,071	3,075	3,080	39/1.3	14.2
Falcon-1024-EMU	Exp SK	26,095	26,101	26,120	109/3.5	120.8
FPU vs EMU	Exp SK	8.50x	<b>8.49x</b>	8.48x	-/-	<b>8.51x</b>

# Benchmarking (Dilithium vs Falcon)



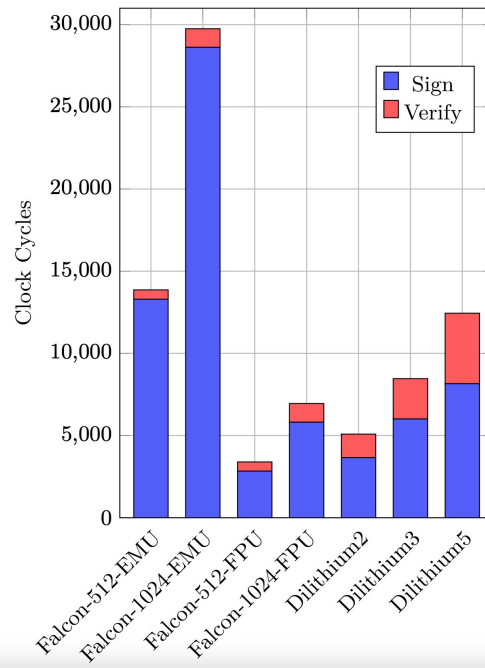
Comparing **Dilithium** and **Falcon** now shows a much different performance profile.



**Falcon-512** now slightly faster than Dilithium2, for both signing and signing+verify runtimes.



**Falcon-1024** also slightly faster than Dilithium5 signing and much faster when combining verify.



# Profiling Falcon (M4 vs M7)

## Performance improvements inside Falcon:

### For key generation:

- iFFT/FFT multiplication 16x improved
- Going from 10m to 0.5m cycles

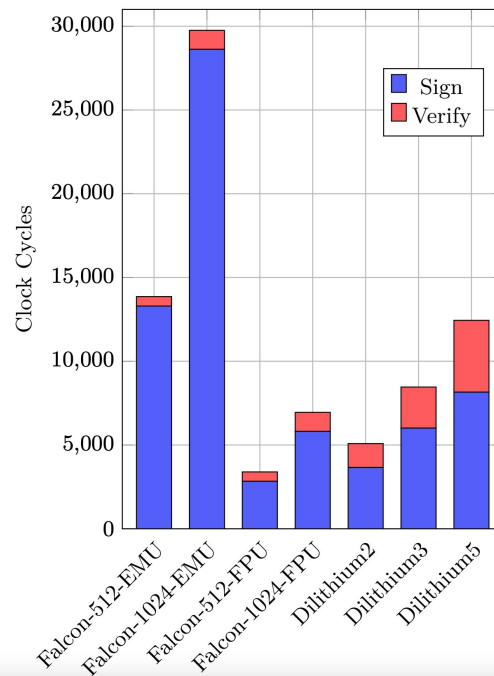
### For both signing modes:

- Fast Fourier sampling >5x improved.
- Going from 16m to <3m cycles.

**Verify times were unchanged.**

**Expand private key improved 12x.**

Going from 11m to <1m cycles.



A large, dark wireframe sphere is positioned on the left side of the frame, partially obscured by the text. The background is a dark, textured space with faint, glowing particles or stars scattered across it.

# 03

## CONSTANT OR ISOCHRONOUS RUNTIME



# Constant-Time Validation

**Floating-point arithmetic is rare in cryptography!**

Thus we thought it was worth looking at...

## We used inline assembly to

- Minimize the unwanted optimizations from the compiler / clobbered registers where necessary.
- This minimizes the effect of surrounding instructions on the operations of interest.
- Which occurred when we tried using C.
- Ensures that all execution is from cache.

This example is for double precision multiplication, i.e., `vmul.f64`, this is repeated for each instruction.

**We tested 4 STM32 development boards.**

```
1 asm volatile (  
2     "vldr d5, %2\n"  
3     "vldr d6, %3\n"  
4     "dmb\n"  
5     "isb\n"  
6     "ldr r1, %1\n"  
7     "vmul.f64 d4, d5, d6\n"  
8     "vmul.f64 d4, d5, d6\n"  
9     "vmul.f64 d4, d5, d6\n"  
10    "vmul.f64 d4, d5, d6\n"  
11    "vmul.f64 d4, d5, d6\n"  
12    "vmul.f64 d4, d5, d6\n"  
13    "vmul.f64 d4, d5, d6\n"  
14    "vmul.f64 d4, d5, d6\n"  
15    "vmul.f64 d4, d5, d6\n"  
16    "vmul.f64 d4, d5, d6\n"  
17    "ldr r2, %1\n"  
18    "subs %0, r2, r1\n"  
19    : "=r"(cycles) : "m"(DWT->CYCCNT),  
20    "m"(r1), "m"(r2) : "r1", "r2",  
21    "d4", "d5", "d6");
```

# Constant-Time Validation

Assembly code uses **two random inputs** for each function.

We found **timing issues** in all double precision FPU instructions across all 4 STM32 boards.

In addition (vadd.f64) runtimes had **16 clocks on avg, standard deviation of 4.1.**

If we generated random values in the same range, such they had the same exponents, the **runtimes were constant and consistent at 10 clock cycles.**

Moreover, when we mixed randomness from two fixed exponent ranges we observed **constant and consistent runtimes of 19 clock cycles.**

```
1 asm volatile (  
2     "vldr d5, %2\n"  
3     "vldr d6, %3\n"  
4     "dmb\n"  
5     "isb\n"  
6     "ldr r1, %1\n"  
7     "vmul.f64 d4, d5, d6\n"  
8     "vmul.f64 d4, d5, d6\n"  
9     "vmul.f64 d4, d5, d6\n"  
10    "vmul.f64 d4, d5, d6\n"  
11    "vmul.f64 d4, d5, d6\n"  
12    "vmul.f64 d4, d5, d6\n"  
13    "vmul.f64 d4, d5, d6\n"  
14    "vmul.f64 d4, d5, d6\n"  
15    "vmul.f64 d4, d5, d6\n"  
16    "vmul.f64 d4, d5, d6\n"  
17    "ldr r2, %1\n"  
18    "subs %0, r2, r1\n"  
19    : "=r"(cycles) : "m"(DWT->CYCCNT),  
20    "m"(r1), "m"(r2) : "r1", "r2",  
21    "d4", "d5", "d6");
```

# Constant-Time Validation

Also tested the ARM Cortex A53 as a previous paper uses Raspberry Pi 3.

Issue found when casting from types **double** to **int64\_t**, op rounds towards zero.



No native instruction to do this on ARMv7.

This can be non-constant time

In LLVM, it isn't, and **leaks the sign**.

**We reported this to the Falcon team and proposed the following fix shown on the right.**

```
1 int64_t cast(double a) {
2     union {
3         double d;
4         uint64_t u;
5         int64_t i;
6     } x;
7     uint64_t mask;
8     uint32_t high, low;
9
10    x.d = a;
11
12    mask = x.i >> 63;
13    x.u &= 0x7fffffffffffffffL;
14
15    // a / 0x1p32f;
16    high = x.d / 4294967296.f;
17
18    // high * 0x1p32f;
19    low = x.d - (double)high * 4294967296.f;
20    x.u = ((int64_t)high << 32) | low;
21
22    return (x.u & ((uint64_t)-1 - mask))
23           | ((-x.u) & mask);
24 }
```

# Takeaways

01

Falcon is super fast on the Cortex M7.

02

Unknown if timing issues can be exploited.

03

Users should consider this thoroughly for each use case.

For example

**Cloudflare currently recommend using Falcon in offline situations.**

```
1 int64_t cast(double a) {  
2     union {  
3         double d;  
4         uint64_t u;  
5         int64_t i;  
6     } x;  
7     uint64_t mask;  
8     uint32_t high, low;  
9  
10    x.d = a;  
11  
12    mask = x.i >> 63;  
13    x.u &= 0x7fffffffffffffffL;  
14  
15    // a / 0x1p32f;  
16    high = x.d / 4294967296.f;  
17  
18    // high * 0x1p32f;  
19    low = x.d - (double)high * 4294967296.f;  
20    x.u = ((int64_t)high << 32) | low;  
21  
22    return (x.u & ((uint64_t)-1 - mask))  
23    | ((-x.u) & mask);  
24 }
```



**THANK  
YOU**

**Paper:**

<https://eprint.iacr.org/2022/405>

**GitHub:**

<https://github.com/jameshoweee/falcon-fpu>

NIST understands that some applications will not work as they are currently designed if the signature and the data being signed cannot fit in a single internet packet. For these applications, the implementation complexity of FALCON's signature generation may not be a concern, but the difficulty of modifying the applications to work with Dilithium's larger signature size may create a barrier to the transition to post-quantum signature schemes. For this reason, NIST decided to standardize FALCON as well. Given FALCON's overall better performance when signature generation does not need to be performed on constrained devices, many applications may prefer to use FALCON over Dilithium, even in cases in which Dilithium's signature size would not be a barrier to implementation.

### Structure:

1. NIST have always considered Falcon to be much slower than Dilithium signing. The ARM Cortex M4 was decided on the benchmarking target for embedded software. Last year we wanted to challenge this, since the M4 does not have the available FPU.
2. Coincidentally and conveniently, the M7 is the next core in the Cortex line, is very similar to the M4 (has a few memory improvements) but has a 64-bit FPU compared to the 32-bit on the M4. So on a board with everything that's needed and with fair benchmarking, which one is faster.
3. Is the statement above by NIST still true.....? No!
4. What parts of the implementations are faster? the FFT stuff etc etc and how much faster vs EMU?



# What's the big deal?

Constant-time and Correctness



**Falcon can heavily use floating point operations for:**

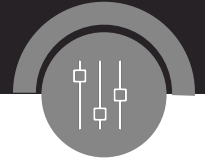
- Additions and subtractions,
- Multiplications, divisions, and accumulation versions,
- Other operations like square root and conversions.



**For determinism, floats should be emulated**



**Native FPUs and code optimizations may yield slight discrepancies**



**Across different platforms and configurations we would like:**

the same message to yield the same signature (for the same secret)