

Crystals-Dilithium 数字签名技术硬件实现综述

崔益军, 李梦雪, 王 隼, 王成华, 刘伟强

(南京航空航天大学集成电路学院, 江苏南京 210016)

摘要: 随着量子计算技术的不断发展, 依赖传统公钥密码体制三大功能(密钥协商/数字签名/公钥加密)的各种应用系统将不再安全. 为应对量子威胁, 以美国国家标准与技术研究院(National Institute of Standards and Technology, NIST)为首的国际标准组织积极征集与部署后量子密码(Post Quantum Cryptography, PQC)算法的标准化工作, 致力于在真正实用型量子计算机问世之前, 提前完成传统公钥密码算法到 PQC 算法的迁移过渡. Crystals-Dilithium 是 NIST-PQC 标准中的一种基于格的数字签名算法, 其安全性高, 运算速度快, 是实现抵抗量子攻击数字签名算法的重要路径之一. 本文从主流 Crystals-Dilithium 数字签名算法的理论基础出发, 从底层关键组件的优化方法和整体硬件构架设计方法着手, 围绕硬件资源优化和性能优化等现有方法和成果对比展开分析介绍, 为研究者们后续研究指明方向, 希望为设计性能与硬件资源均衡的后量子数字签名密码芯片提供有力参考.

关键词: 后量子密码; 格密码; Crystals-Dilithium 数字签名; 硬件实现; 优化方案

基金项目: 国家自然科学基金(No.62134002, No.62104107, No.62371226)

中图分类号: TP393 **文献标识码:** A **文章编号:** 0372-2112(2025)07-2558-21

电子学报 URL: <http://www.ejournal.org.cn> **DOI:** 10.12263/DZXB.20240505

A Survey of Hardware Implementation for Crystals-Dilithium Digital Signature in the Post Quantum Cryptography

CUI Yi-jun, LI Meng-xue, WANG Bei, WANG Cheng-hua, LIU Wei-qiang

(College of Integrated Circuits, Nanjing University of Aeronautics and Astronautics, Nanjing, Jiangsu 210016, China)

Abstract: With the continuous development of quantum computing technology, various application systems relying on the three major functions of traditional public key cryptography (key agreement/digital signatures/public key encryption) will no longer be secure. In response to the quantum threat, international standardization organizations led by the United States NIST (National Institute of Standards and Technology) are actively soliciting and deploying the standardization work of post-quantum cryptography (PQC) algorithms, aiming to complete the migration from traditional public key cryptography algorithms to PQC algorithms before truly practical quantum computers emerge. Crystals-Dilithium is one of the lattice-based digital signature algorithms in the NIST-PQC standard, which features high security and fast computation speed, making it an important path to implement digital signature algorithms resistant to quantum attacks. This paper commences with the theoretical foundations of the mainstream Crystals-Dilithium digital signature algorithm. It delves into optimization methods for underlying key components and overall hardware architecture design, focusing on hardware resource optimization and performance enhancement. The paper contrasts and analyzes existing methods and outcomes, aiming to clarify the direction for subsequent research. It aspires to provide a robust reference for the design of post-quantum digital signature cryptographic chips that balance performance with hardware resource efficiency.

Key words: post quantum cryptography; lattice-based cryptography; Crystals-Dilithium digital signature; hardware implementation; optimization scheme

Foundation Item(s): National Natural Science Foundation of China (No.62134002, No.62104107, No.62371226)

1 引言

当今信息安全领域广泛使用的公钥密码体制主要是基于难以求解的数学困难问题而构造的。例如: Rivest-Shamir-Adleman^[1]是基于大整数分解问题, Diffie-Hellman^[2]和 ElGamal^[3]是基于有限域上离散对数问题, 椭圆曲线密码(Elliptic Curve Cryptography, ECC)算法^[4]是基于椭圆曲线离散对数问题。在经典的计算模型下, 这些底层数学困难问题的攻击复杂度至少是亚指数级的^[5]。1994年, 美国数学家 Peter Shor 提出 Shor 算法, 能够在多项式时间内求解大整数分解问题和离散对数问题^[6]。这表明一旦真正实用型量子计算机问世, 现有传统公钥密码体制的安全性将面临坍塌的风险。为应对量子计算机对传统公钥密码体制带来的挑战, 2016年12月, 美国国家标准技术研究院(National Institute of Standards and Technology, NIST)开展后量子密码学(Post-Quantum Cryptography, PQC)算法标准征集工作^[7], 并计划通过几轮的标准化竞争过程评选出能够抵抗量子攻击的公钥密码体制。NIST的标准化评定对于算法的评估主要涵盖三个方向: 一是评估候选密码算法的安全性, 包括算法的抗攻击性、密码强度以及对各种安全威胁的防护能力; 二是评估算法的运行效率, 包括密钥生成速度、封装/解封装速度、签名/验签速度、内存消耗等指标; 三是评估算法的实际应用性, 包括算法的实现难度、对硬件和软件的适应性, 以及算法的易用性和可扩展性。经过多轮评审, NIST于2022年7月公布了被采纳使用且待标准化的PQC方案^①, 宣布 Crystals-Kyber、Crystals-Dilithium、FALCON 和 SPHINCS+ 这4种PQC算法将进入待标准化进程中, 其中基于格的PQC方案有 Crystals-Kyber 和 Crystals-Dilithium, 前者是实现密钥封装机制, 而后者是实现数字签名协议。2023年8月, NIST发布3份后量子密码算法标准草案, 将 Crystals-Kyber 优化并更名为 ML-KEM, Crystals-Dilithium 优化并更名为 ML-DSA (Modular Lattice-based Digital Signature Algorithm, ML-DSA), SPHINCS+ 优化并更名为 SLH-DSA (StateLess Hash-based Digital Signature Algorithm)。2024年8月 NIST 正式确定此3个算法为标准后量子密码算法方案。2024年11月, NIST 发布了过渡到后量子密码学标准的初始公开草案《Transition to Post-Quantum Cryptography Standards》^②, 包含迁移的路线与时间表。NIST 希望到2035年将政府机构的加密系统替换为后量子密码体制。

后量子密码方案所基于的数学原理主要包括以下五种: 基于格(Lattice-based)^[8]、基于哈希(Hash-based)^[9]、基于编码(Code-based)^[10]、基于多变量(Multivariate-based)^[11]和基于超奇异同源(Isogeny-based)^[12]。由于不同的底层数学结构和设计门限方案,

这些密码方案具有差异化的参数和性能。表1从底层数学结构、优点、缺点和发展情况这四个方面介绍了主流后量子密码算法, 可以看出基于格的PQC算法相较于其他算法具有更好的表现, 在安全性、公私钥大小、运算速度上可达到较好的平衡, 是目前最有应用前景的PQC算法之一。

在PQC征集过程中, 一些PQC算法不断遭受挑战。如2022年2月, 晋级NIST的PQC决赛第三轮的Rainbow算法在其某个参数设置下被笔记本电脑成功破解^[13]; 同年4月, 3款基于LWE(Learning With Errors)的NIST后量子密码算法被曝再遭挑战^[14]; 2022年7月31日, 比利时鲁汶大学的研究人员^[15]通过Magma计算软件, 彻底破解了基于超奇异同源的密钥协商算法超奇异同源密钥封装(Supersingular Isogeny Key Encapsulation, SIKE)。截至目前, 基于格的后量子密码算法依然被国内外研究学者和企业高度认可。近些年, 针对格密码算法(Crystals-Kyber 和 Crystals-Dilithium)软硬件实现的研究与应用越来越多, 随之也面临迁移过程中的软硬件灵敏性和适配性等问题。

2021年12月, 全球领先的半导体企业进阶精简指令集机器(Advanced RISC Machines, ARM)推出了关于后量子密码技术的白皮书^③, 从中可以看到国际芯片企业对未来网络空间安全基础技术的前瞻性布局。事实上, 包括Intel等国际芯片企业在内, 均早已在关注未来芯片的抗量子计算的攻击。2022年4月, 美国众议院推动联邦政府信息技术系统向PQC算法过渡, 我国也正在积极推动相关标准化制定和产业应用。随着NIST对PQC算法标准化的推进, 关于PQC算法的硬件实现研究也步入了快速发展的轨道, 但是从当今的公钥密码体制向PQC算法迁移是一项既耗时又具有挑战性的工作。最近的一些报道也表明, 一些领先的科技公司正在加大对硬件PQC实现的研发和投入力度, 这一举措被视为未来网络安全战略的重要组成部分。因此, 可以预见, 硬件实现将成为未来PQC迁移的关键推动力量, 为各行业提供更为高效、更为安全的解决方案, 从而保障了数据的隐私和安全。目前, 在PQC算法的硬件实现方面相关研究较少。在算法的硬件设计过程中, 仍然需要在算法及计算架构层面, 解决性能资源、速度和功耗不均衡, 兼容性和安全性低等一系列的硬件实现问题。

考虑到PQC在从云到端的部署与应用中都离不开PQC芯片的支撑, 因此设计出实用性、灵活性、高效性

① <https://www.nist.gov/publications/status-report-second-round-nist-post-quantum-cryptography-standardization-process>

② <https://csrc.nist.gov/pubs/ir/8547/ipd>

③ <https://community.arm.com/arm-research/m/resources/1002>

表 1 主流后量子密码算法

类别	安全性依赖的数学问题	优点	缺点	发展情况
基于格	依赖于格中困难问题	公私钥尺寸小、计算速度快,且能被用于构造多种密码学原语,因此更适用于实际应用环境	密钥封装的交互过程中,传输的密文尺寸较大	格密码算法被广泛认为是PQC算法征集中的主力军,目前已有两个格密码算法被标准化,未来Falcon也在标准化进程中
基于哈希	依赖于哈希算法的一些安全性质,例如单向性(抗原像攻击)、弱抗碰撞性(抗第二原像攻击)和伪随机性等	公钥很小、安全假设少、安全性强	签名尺寸较长,功能上只能构造签名	基于哈希的密码算法更适用于数字签名,但底层数学结构少,目前Sphincs+已被标准化
基于编码	使用错误纠正码对加入的随机性错误进行纠正和计算,对于随机线性码的解码问题是非确定性多项式时间(Non-deterministic Polynomial time, NP)难	底层数学困难问题研究时间长,安全性更经得起考验	公钥大、密钥生成慢	后量子密码中相对具有竞争力的密码算法,但在实用化方面有待优化提升
基于多变量	求解有限域上非线性方程组 NP 难问题	签名验签速度快、消耗资源少	公钥大	NIST-PQC的附加数字签名方案第二轮中提交了多个基于多变量的密码算法,在多变量赛道大概率有算法被标准化
基于超奇异同源	寻找任意两条椭圆曲线之间的同源问题	公私钥或签名验签的尺寸较小	计算速度慢、代数结构复杂、缺乏可证明安全性	SIKE被攻破,不过交换超奇异同源迪菲-赫尔曼(Commutative Supersingular Isogeny Diffie-Hellman, CSIDH)和超奇异二次同源签名(Supersingular Quadratic Isogeny signatures, SQIsign)算法仍未被攻破,且SQIsign已进入NIST-PQC的附加数字签名方案第二轮

与安全性有机统一的PQC芯片至关重要.在当前网络信息互联时代,庞大的电子商务体系、银行和金融服务、电子数字签名与网址安全访问等应用广泛,其底层大量使用数字签名算法.未来PQC数字签名体制必将应用到各行各业,而其签名算法硬件实现的资源和性能是衡量PQC数字签名算法应用落地的重要指标.本文从主流Crystals-Dilithium数字签名算法的理论基础出发,从底层关键组件的优化方法和整体硬件构架设计方法着手,围绕硬件资源优化和性能优化等现有方法和成果对比展开综述,为数字签名芯片的设计实现与应用提供参考.

2 背景介绍

本章节将围绕Crystals-Dilithium数字签名算法的理论展开介绍.首先,介绍本章节算法中出现的符号,接着介绍算法的详细实现流程,最后对算法中的环多项式乘法算法进行总结分析,为后续章节中Crystals-Dilithium数字签名算法的高效硬件实现与关键技术研究提供理论依据.

2.1 符号说明

本节主要用到的符号说明:设 n 和 q 是两个整数, n 表示为2的幂, q 表示为满足 $q \equiv 1 \pmod{2n}$ 的质数.设 \mathbb{Z}

表示整数的集合, R 表示环多项式环 $\mathbb{Z}[x]/(x^n+1)$,其中 (x^n+1) 是 $\mathbb{Z}[x]$ 的理想.另外 \mathbb{Z}_q 定义为 $\mathbb{Z}/q\mathbb{Z}$, R_q 是由 $\mathbb{Z}_q[x]/(x^n+1)$ 得到的商环. R (或 R_q)中的元素为多项式,为了便于区分,这里用斜体化的大写字母表示,例如 F ,相反,向量用粗体的、斜体化的小写字母表示,如 \mathbf{x} .每个多项式元素,无论是 $F \in R$ 还是 $F \in R_q$,都可以只表示为 $F = \sum_{i=0}^{n-1} F_i X^i$,其中 F_i 属于 $i \in [0, n)$ 的 \mathbb{Z} (或 \mathbb{Z}_q).

R (或 R_q)中的矩阵用粗体的、斜体化的大写字母表示,例如 A .对于多项式 $F \in R$,它的无穷范数可以表示为 $\|F\|_\infty = \max\{|F_i|\}$,同样的对于一个向量 $\mathbf{x} = (x_0, x_1, \dots, x_{n-1}) \in R^n$,它的无穷范数可以表示为 $\|\mathbf{x}\|_\infty = \max_{0 \leq i \leq n-1} |x_i|$.

2.2 Crystals-Dilithium 数字签名算法

Crystals-Dilithium是一种基于格的数字签名算法,算法采用在环多项式 $R_q = \mathbb{Z}_q[x]/(x^{256}+1)$ 上的Fiat-Shamir范式构造^[16].其安全性基于模误差学习问题(Module-Learning With Errors, M-LWE)和模短整数解(Module-Short Integer Solution, M-SIS)问题^[17].2023年8月,NIST发布后量子密码学标准草案,Crystals-Dilithium方案被修改和重新命名为ML-DSA,但其算法的整体架构与流程未发生改变,仅为了确保ML-DSA满

足特定的安全属性,如不可伪造性和非否认性,针对 Crystals-Dilithium 的某些参数进行了调整,例如增加了公钥 tr 的长度到 512 位,以及增加了 \tilde{c} 的长度到 384 位和 512 位等. 为了保持算法的最新性,本文介绍 ML-DSA 的算法流程,包括密钥生成、签名及验签环节.

表 2 展示了三组 ML-DSA 参数集的算法参数. 每个参数集都用于 ML-DSA 算法的密钥生成、签名和验证的所有参数分配值. 表中第 1 行的参数集 ML-DSA-44 属于安全等级 2, ML-DSA-65 属于安全等级 3, 而 ML-DSA-87 属于安全等级 5. 三种不同的参数集中模数(q)、采样系数(d)相同. 其中, (k, l) 代表着多项式矩阵 A 的维度, 通过改变 (k, l) 的值, 可以调整签名的安全等级. η 代表着私钥的范围, β 是验证 $if c = H(\rho, t_1, w_1, \mu)$ and $\|z\|_\infty < \gamma_1 - \beta$ 环节中的范数误差, ω 代表着提示 h 中的最大值, Repetitions 代表着签名算法主循环的预期重复次数. 从表中可以看出, 随着安全等级不断上升, 范数误差的值在逐渐减小, 签名信息也会更加准确.

表 2 ML-DSA 数字签名算法参数列表

序号	参数	ML-DSA-44	ML-DSA-65	ML-DSA-87
1	Claimed Security Level	2	3	5
2	(k, l)	(4,4)	(6,5)	(8,7)
3	q	8 380 417	8 380 417	8 380 417
4	d	13	13	13
5	τ	39	49	60
6	λ	128	192	256
7	γ_1	217	219	219
8	γ_2	($q-1$)/88	($q-1$)/32	($q-1$)/32
9	η	2	4	2
10	$\beta = \tau \cdot \eta$	78	196	120
11	ω	80	55	75
12	$\log\left(\frac{256}{\tau}\right) + \tau$	192	225	257
13	Repetitions	4.25	5.1	3.85

密钥生成过程. 密钥生成过程主要完成字节字符串公钥 pk 和私钥 sk 的生成, 如算法 1 所示. 首先, 算法生成一个 256 的随机种子 ξ , 通过 SHAKE-256 哈希函数来产生其他随机值 ρ, ρ', K . 其中, 公共的随机 256 位种子 ρ 经过一个伪随机采样函数 ExpandA 将其展开到多项式矩阵 $A \in R_q^{k \times l}$ 中并通过后续计算转化为数论变换 (Number Theoretic Transform, NTT) 的表达式 \hat{A} , 其中 \hat{A} 矩阵的项都是 R_q 环中的多项式. 接着通过伪随机采样函数 ExpandS 展开另一个私有的随机 256 位种子 ρ' , 产生满足在 $[-\eta, \eta]$ 之间具有均匀随机系数的密钥向量 s_1 和 s_2 . 紧接着, 通过 NTT 矩阵多项式计算操作得到公共

值 $t = As_1 + s_2$, 再通过 Power2Round 函数将其分解为高低位, 其中 t 的低阶比特可以从少量签名中重建, 因此为了提升运算的性能, 算法将每个系数中 d 个最不重要的低阶比特丢弃从而产生压缩多项式系数向量 t_1 . 最后算法通过字符串编码操作译为公钥 pk 和私钥 sk .

算法 1 ML-DSA.KeyGen()

输出: 公钥 $pk \in B$, 私钥 $sk \in B$

1. $\xi \leftarrow \{0, 1\}^{256}$
2. $(\rho, \rho', K) \in \{0, 1\}^{256} \times \{0, 1\}^{512} \times \{0, 1\}^{256} \leftarrow H(\xi, 1024)$
3. $\hat{A} \leftarrow \text{ExpandA}(\rho)$
4. $(s_1, s_2) \leftarrow \text{ExpandS}(\rho')$
5. $t \leftarrow \text{NTT}^{-1}(\hat{A} \circ \text{NTT}(s_1)) + s_2$
6. $(t_1, t_0) \leftarrow \text{Power2Round}(t, d)S$
7. $pk \leftarrow \text{pkEncode}(\rho, t_1)$
8. $tr \leftarrow H(\text{BytesToBits}(pk), 512)$
9. $sk \leftarrow \text{skEncode}(\rho, K, tr, s_1, s_2, t_0)$
10. RETURN (pk, sk)

签名过程. 签名过程将被编码为字节串的私钥 sk 和位字符串消息 M 作为输入, 为消息 M 生成一个签名 σ , 如算法 2 所示. 首先签名者从私钥 sk 中通过 skDecode 函数提取所需要的 $\rho, K, tr, s_1, s_2, T_0$ 信息, 其中 ρ 是 256 位的公共随机种子; K 是 256 位的私有随机种子, tr 是 512 位的公钥的哈希值; s_1, s_2 是私钥的多项式向量; t_0 是多项式向量 t 的每个系数 d 的最小有效位多项式向量. 接着, 签名者将 s_1, s_2, t_0 通过计算转化为 NTT 的表达式 $\hat{s}_1, \hat{s}_2, \hat{t}_0$, 展开 ρ 到与密钥生成中相同的矩阵 $A \in R_q^{k \times l}$ 并转化为 NTT 的表达式 \hat{A} . 在消息 M 被签名之前, 签名者将它与公钥 tr 哈希连接, 并使用哈希函数 H 生成 512 位消息代表 μ , 接着计算私有随机种子 ρ' . 签名算法的主要部分为生成签名的拒绝采样循环, 如果签名可能泄露私钥信息, 则判断为无效签名, 算法将重复循环, 直到产生有效签名, 然后将其编码为字节字符串并输出. 拒绝采样循环所涉及的主要步骤有: 首先, 算法签名者使用 ExpandMask 函数、种子 ρ' 和一个计数器 κ , 从多项式向量的子集中伪随机地采样系数在 $[-\gamma_2, \gamma_2]$ 范围内的多项式向量 $y \in R_q^l$. 在获得签名者样本 y 后, 签名者首先计算 $w = Ay$, 再根据公式 $w = 2\gamma_2 \cdot w_1 + w_0$, 通过使用 HighBits 函数, 四舍五入到最接近的 $2\gamma_2$ 的倍数来计算承诺 w_1 ; 在计算完 w_1 后, 连接消息 μ 和编码过后的 w_1 值, 利用哈希函数创建其散列值并转化哈希形式 \tilde{c} , 再用 \tilde{c}_1 表示 \tilde{c} 的前 256 位. 比特串 \tilde{c}_1 被用来伪随机地采样一个多项式 $c \in R_q$, 其系数为 $-1, 0$ 或 1 , 且非零系数的数量为 τ . 在得到散列值 c 后, 计算签名 $z = y + cs_1$, 如果此时直接输出 z , 将可能导致密钥被泄

露,那么签名方案将是不安全的.因此,为了避免 z 会泄露密钥的信息,算法使用拒绝抽样的方式检验前面是否有效.如果检查通过,签名者可以计算一个提示多项式 h ,这将允许验证者使用压缩的公钥(以及签名的其他组件)重构 w_1 .最后签名者输出最终的签名,即承诺 \tilde{c} 、响应 z 和提示 h 的字节编码 σ .

验签过程.验签过程将被编码为字节字符串的公钥 pk 、位字符串的消息 M 以及字节字符串的签名 σ 作

算法2 ML-DSA.Sign(sk, M)

输入: 私钥 $sk \in B$, 消息 $M \in \{0, 1\}^*$
 输出: 签名 $\sigma \in B$

1. $(p, K, tr, s_1, s_2, t_0) \leftarrow \text{skDecode}(sk)$
2. $\hat{s}_1 \leftarrow \text{NTT}(s_1)$
3. $\hat{s}_2 \leftarrow \text{NTT}(s_2)$
4. $\hat{t}_0 \leftarrow \text{NTT}(t_0)$
5. $\hat{A} \leftarrow \text{ExpandA}(\rho)$
6. $\mu \leftarrow H(tr \parallel M, 512)$
7. $rnd \leftarrow \{0, 1\}^{256}$
8. $\rho' \leftarrow H(K \parallel rnd \parallel \mu, 512)$
9. $\kappa \leftarrow 0$
10. $(z, h) \leftarrow \perp$
11. WHILE $(z, h) = \perp$ DO
12. $y \leftarrow \text{ExpandMask}(\rho', \kappa)$
13. $w \leftarrow \text{NTT}^{-1}(\hat{A} \circ \text{NTT}(y))$
14. $w_1 \leftarrow \text{HighBits}(w)$
15. $\tilde{c} \in \{0, 1\}^{2\lambda} \leftarrow H(\mu \parallel w1\text{Encode}(w_1), 2\lambda)$
16. $(\tilde{c}_1, \tilde{c}_2) \in \{0, 1\}^{256} \times \{0, 1\}^{2\lambda-256} \leftarrow \tilde{c}s$
17. $c \leftarrow \text{SampleInBall}(\tilde{c}_1)$
18. $\hat{c} \leftarrow \text{NTT}(c)$
19. $\langle\langle cs_1 \rangle\rangle \leftarrow \text{NTT}^{-1}(\hat{c} \circ \hat{s}_1)$
20. $\langle\langle cs_2 \rangle\rangle \leftarrow \text{NTT}^{-1}(\hat{c} \circ \hat{s}_2)$
21. $z \leftarrow y + \langle\langle cs_1 \rangle\rangle$
22. $r_0 \leftarrow \text{LowBits}(w - \langle\langle cs_2 \rangle\rangle)$
23. IF $\|z\|_\infty \geq \gamma_1 - \beta$ or $\|r_0\|_\infty \geq \gamma_2 - \beta$ THEN $(z, h) \leftarrow \perp$
24. ELSE
25. $\langle\langle ct_0 \rangle\rangle \leftarrow \text{NTT}^{-1}(\hat{c} \circ \hat{t}_0)$
26. $h \leftarrow \text{MakeHint}(-\langle\langle ct_0 \rangle\rangle, w - \langle\langle cs_2 \rangle\rangle + \langle\langle ct_0 \rangle\rangle)$
27. IF $\|ct_0\|_\infty \geq \gamma_2$ or the number of 1's in h is greater than ω
28. THEN $(z, h) \leftarrow \perp$
29. END IF
30. END IF
31. $\kappa \leftarrow \kappa + \ell$
32. END WHILE
33. $\sigma \leftarrow \text{sigEncode}(\tilde{c}, z \bmod^+ q, h)$
34. RETURN σ

为输入,验证消息 M 的签名 σ 是否正确,如算法3所示.验证者首先从公钥 pk 中提取公共随机种子 ρ 和压缩多项式向量 t_1 ;接着从签名 σ 中提取签名者的承诺 \tilde{c} 、响应 z 和提示 h .如果验证者发现提示 h 未被正确地字节编码,则用符号“ \perp ”表示,验证算法立即返回false,表示签名无效.假设签名成功地从其字节编码中提取出来,验证者通过伪随机方法从公共随机种子 ρ 中拓展多项式矩阵 \hat{A} ,并通过其对 pk 的哈希结果 tr 和消息 pk 的串联进行哈希运算来创建一个消息代表 μ .然后,验证者尝试根据公钥 pk 和签名 σ 重构签名者的多项式向量承诺 w_1 ,其中重构的值被称为 w'_1 .最后,验证者检查签名者的响应 z 和签名者的提示 h 是否有效以及重构的 w'_1 与签名者的承诺哈希 \tilde{c} 是否一致,即检查:(1)所有系数是否足够小(即在范围 $[-\gamma_2, \gamma_2]$ 内);(2) h 是否包含不超过 w 个非零系数;(3) \tilde{c} 是否与消息代表 μ 串联和表示为位字符串的 w'_1 的哈希 \tilde{c}' 匹配.如果所有这些检查都成功,则返回true,否则返回false.

算法3 ML-DSA.Verify(pk, M, σ)

输入: 私钥 $sk \in B$, 消息 $M \in \{0, 1\}^*$, 签名 $\sigma \in B$
 输出: Boolean.

1. $(\rho, t_1) \leftarrow \text{pkDecode}(pk)$
2. $(\tilde{c}, z, h) \leftarrow \text{sigDecode}(\sigma)$
3. IF $h = \perp$ THEN RETURN FALSE
4. END IF
5. $\hat{A} \leftarrow \text{ExpandA}(\rho)$
6. $tr \leftarrow H(\text{BytesToBits}(pk), 512)$
7. $\mu \leftarrow H(tr \parallel M, 512)$
8. $(\tilde{c}_1, \tilde{c}_2) \in \{0, 1\}^{256} \times \{0, 1\}^{2\lambda-256} \leftarrow \tilde{c}$
9. $c \leftarrow \text{SampleInBall}(\tilde{c}_1)$
10. $W'_{\text{APPROX}} \leftarrow \text{NTT}^{-1}(\hat{A} \circ \text{NTT}(z) - \text{NTT}(c) \circ \text{NTT}(t_1 \cdot 2^d))$
11. $W'_1 \leftarrow \text{UseHint}(h, W'_{\text{APPROX}})$
12. $\tilde{c}' \leftarrow H(\mu \parallel w1\text{Encode}(W'_1), 2\lambda)$
13. RETURN $\left[\|z\|_\infty < \gamma_1 - \beta \right] \text{AND} \left[\tilde{c} = \tilde{c}' \right] \text{AND} \left[\text{number of 1's in } h \text{ is } \leq w \right]$

表3展示了每个参数集相对应的密钥和签名的大小.从此表中可以看到随着安全等级的不断上升,密钥和签名的规模也在逐渐增大.在存储ML-DSA公钥和私钥时还可以采取不同的优化方式.如果有较多的硬件存储资源,可以预先计算并存储多项式矩阵 \hat{A} ,以加快签名和验证的速度.如果希望减少私钥所需的存储空间,则只需要存储32字节的种子 ζ ,在签名生成过程中计算私钥的其他部分.

2.3 Crystals-Dilithium 环多项式乘法算法

Crystals-Dilithium 的多项式乘法在环多项式 R_q 上计算,值得注意的是:(1)多项式的最高项次数必须小于 n ;(2)多项式的系数必须在模 q 范围内.若乘法的结

表 3 ML-DSA 数字签名密钥和签名大小 单位:Byte

参数集	私钥	公钥	签名
ML-DSA-44	2 528	1 312	2 420
ML-DSA-65	4 000	1 952	3 293
ML-DSA-87	4 864	2 592	4 595

果就在多项式环上则不需要进行处理,否则需要对结果进行模操作,从而映射到多项式环中。

给定两个多项式 $A = \sum_{i=0}^{n-1} A_i x^i, B = \sum_{j=0}^{n-1} B_j x^j$, 若用公式法计算多项式乘法,如式(1)所示:

$$C(x) = A(x) \times B(x) = \left[\sum_{i=0}^{n-1} \sum_{j=0}^{n-1} A_i B_j x^{i+j} \right] \bmod (x^n + 1) \quad (1)$$

$$= \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} (-1)^{\lfloor (i+j)/n \rfloor} A_i B_j x^{(i+j) \bmod n}$$

该方式也被称为教科书(Schoolbook)算法,需要 $O(n^2)$ 次乘法。在硬件实现时,如果采用单个乘法单元实现占用资源较少,但运算耗时较长^[18];如果采用并行结构,虽然能够节省运算时间,但是会导致过多的资源消耗^[19]。

NTT 算法需要 $O(n \log n)$ 次乘法,相比较传统的计算方式更加快速和高效,灵活性也更高^[20]。NTT 算法源于快速傅里叶变换(Fast Fourier Transform, FFT), NTT 去除了 FFT 中复数部分复杂的浮点数和乘加计算,所有运算都在有限域上,且用整数原根代替了 FFT 中的复数形式的单位根(可统称为旋转因子),除了对模数有一定限制,NTT 能够解决大多数格密码算法中的多项式乘法。NTT 变换规则如式(2)所示:

$$\hat{A}_i = \sum_{j=0}^{n-1} A_j w^{ij} \bmod q, i=0, 1, \dots, n-1 \quad (2)$$

将得到的 NTT 域上的数据再返回多项式环上称为 NTT 逆变换(Inverse Number Theoretic Transform, INTT),其变换规则如式(3)所示:

$$A_i = n^{-1} \sum_{j=0}^{n-1} \hat{A}_j w^{-ij} \bmod q, i=0, 1, \dots, n-1. \quad (3)$$

利用 NTT 计算有限域上的多项式乘法 $C(x) = A(x) \times B(x)$ 时,首先对 $A(x)$ 和 $B(x)$ 利用式(2)进行 NTT 变换获得 NTT 域上的表达式 $\hat{A}(x)$ 和 $\hat{B}(x)$, 然后通过 $\hat{C}(x) = \hat{A}(x) \circ \hat{B}(x)$ 将乘法简化为逐点乘法,最后利用式(3)将结果 $\hat{C}(x)$ 恢复到多项式环运算内。整个过程可以用以下公式概括:

$$A(x) \times B(x) = \text{NTT}^{-1}(\text{NTT}(A) \circ \text{NTT}(B)) \quad (4)$$

除了常用的 NTT 算法和 Schoolbook 算法,解决多项式乘法还可以使用 Karatsuba 算法^[21], Karatsuba 算法通过将较长的多项式拆分成多个较短的多项式,并重复

利用较短多项式计算结果,减少了乘法次数,降低了计算复杂度,但在硬件实现时会带来较大的资源开销,影响系统的整体性能。NTT 算法是目前所有环多项式乘法算法中计算复杂度最低的一类算法,并可将 FFT 的硬件结构拓展到 NTT 上,因此硬件结构灵活多变,可以根据不同的应用场景选择相应的 NTT 硬件结构,使得资源消耗和性能达到平衡。

3 Crystals-Dilithium 底层关键组件的优化方案

针对 Crystals-Dilithium 数字签名不同的安全参数和不同的实现平台,例如现场可编程逻辑门阵列(Field Programmable Gate Array, FPGA)平台和专用集成电路(Application Specific Integrated Circuit, ASIC)平台,国内外已经有一些成果。在 Crystals-Dilithium 算法中,哈希和伪随机数扩展、多项式乘法占据了整个算法超过 90% 的资源和时间消耗^[22]。为了提高硬件实现中的性能和效率,对这些底层关键组件进行优化是至关重要的。本章节将会对这些关键组件的优化进行详细阐述和对比分析。

3.1 模乘的优化方案

多项式乘法占据了 Crystals-Dilithium 算法实现近 50% 的运行时间^[23]。NTT 和 INTT 作为多项式乘法的核心加速算法,其中模乘操作占据了主要的计算比重。因此,对模乘操作的优化对于提升 NTT 实现的性能至关重要。其中,优化模乘操作不仅涉及硬件实现结构的优化和资源的权衡,也包括算法层面的改进,通过改进多项式乘法模块中的模约简算法,主要的算法有 Montgomery 模乘算法^[24]、Karatsuba 分治算法^[22]、Barrett 模约简算法^[25]和 Plantard 模约简算法^[26],从而提升算法的执行效率和硬件实现的性能。

Land 等人^[27]采取了一种经典的模约简策略,利用 $2^{23} \equiv 2^{13} - 1 \bmod q$ 的关系对模数进行约简,将一个 23 bit 的数 $r[22:0]$ 利用模数的特性约简为 $r[22:19] \cdot 2^9 + r[18:0]$ 。作为首个针对 Xilinx Artix-7 平台提出的设计, Land 等人^[27]通过使用大量数字信号处理器(Digital Signal Processors, DSPs)资源对模约简的结构进行优化设计,获得最快计算速度的优化方案。具体来说,为了增加计算的吞吐量,作者巧妙地利用移位寄存器实例来延迟第三个 DSP 的输入,如图 1 所示。这种设计在模乘法阶段仅使用了 524 个查找表(Look-Up-Tables, LUTs)以及 759 个触发器(Flip-Flops, FFs),但是在单个模乘法单元中仅使用了 4 个 DSPs。这样的做法虽然可以让整个模乘达到较高的工作频率,但是同时大量 DSPs 的使用,也会导致资源消耗的显著增加。Gupta 等人^[28]也参考这样的递归利用模数的性质进行约简的方式,不

同的是,为了降低资源的消耗,Gupta 等人^[28]利用 LUTs 和 FFs 代替了其中 DSPs 的使用. 这种替代不仅

减少了资源的占用,而且为低功耗和低成本的硬件实现提供了一种有效的思路.

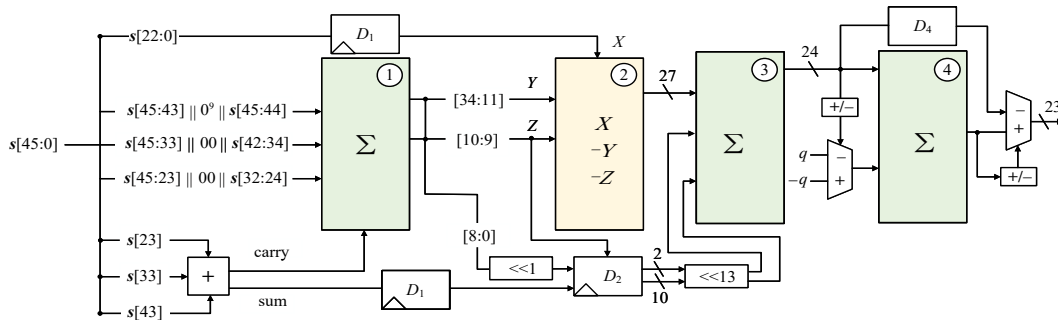


图1 基于数个 DSP 的模约简单元^[27]

Zhao 等人^[22]则在该约简方法的基础上针对 Crystals-Dilithium 算法的三个模数进行了深入研究,并针对算法中涉及的模数提出了进一步的优化策略. 首先,他们将模数转换为一个标准的有符号数字(Canonical Signed Digit, CSD)表示法^[29],这一步骤利用了模数的固有属性来压缩需要约简的数的位宽. 通过重复压缩过程,直到结果小于模量的两倍,从而显著减少了数据的位宽需求. 在压缩步骤完成后,通过执行条件减法操作来获得最终的约简结果,如图2所示.

在模乘的算法实现中,除了从硬件结构和资源优化方面对模约简过程做优化之外,众多研究也从算

法角度入手,致力于改进多项式乘法模块中的模约简算法,以减少所需的时钟周期并提升乘法计算的性. Li 等人^[25]将 Karatsuba 中的分治算法应用在 Crystals-Dilithium 的大数乘法中,实现了显著的性能提升. 由于 DSP 单元无法一次性处理 23 比特乘法, Li 等人^[25]采用分治设计可以使所有 DSPs 单元的计算在一个时钟周期中完成,极大降低了关键路径中的乘法复杂度的同时,显著提高了模块并行度,如图3所示. 此外,该模乘法器架构中部署了三级流水线的设计,其中第一级完成 Karatsuba 算法的中间结果计算,第二级进行移位和加减运算并使用寄存器暂存中间结果,打断关键路径. 第三级完成最终的模约简. 采用该流水线的设计使得各级延迟均衡,吞吐率达到最大化. 整体电路在满足时序约束的同时实现了面积优化. 这种优化策略不仅减少了所需的时钟周期,而且通过并行处理提高了整体的计算吞吐量,但是这样的处理方式不可避免地造成了较大资源的消耗.

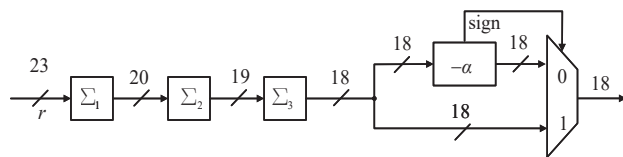


图2 模数为 8380417 的模约简单元^[22]

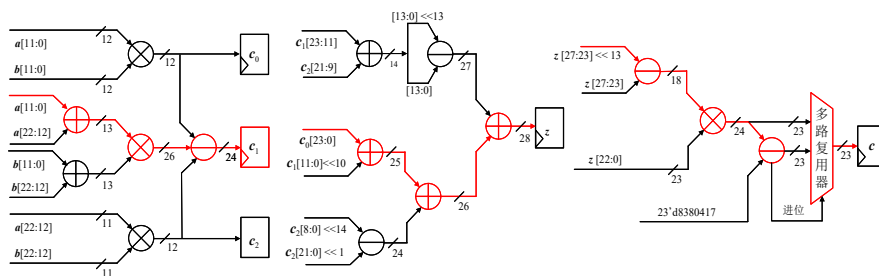


图3 分治法模乘法器架构^[25]

Hwang 等人^[30]归纳了整数近似的概念并且利用 Barrett 算法的近似性质在 Cortex-M3 和 8 位 AVR 上实现 Crystals-Dilithium 的模乘. 在标准 Barrett 乘法中,为了计算 $ab \bmod q$,通常先计算一个近似商 $\tilde{q} = \text{floor}\left(\frac{abR}{q}\right)$,其中 R 是 2 的幂,使得 $ab - q\tilde{q}$ 落在区间 $[-q, q]$ 内. 结合商的近似特性, Hwang 等人^[30]提出可以通过放宽边界条

件,即让 $ab - q\tilde{q}$ 落在一个更宽的区间 $[-B, B]$ 内,其中 $q \leq B \leq R$,降低商 \tilde{q} 的计算成本,因此作者在计算 \tilde{q} 时直接抛弃了部分低位,从而简化了高位乘法的模拟,用更少的乘法和移位操作来近似模拟高位乘法. 此外,作者证明了在 Cortex-M3 和 AVR 平台上,当乘法高指令和乘法长指令不可用或速度较慢时,只需 1 次高位乘法和 2 次低位乘法的 Barrett 乘法比需要 2 次高位乘法或长乘法

和 1 次低位乘法的 Montgomery 乘法更有优势。

在 Cortex-M3 上,基于 Barrett 乘法的 NTT/INTT 实现相比最新的基于 Montgomery 乘法的实现,常数时间版本快 1.38~1.51 倍,变量时间版本快 1.10~1.21 倍.在 8 位的 AVR 平台上,基于 Barrett 乘法的 NTT/INTT 实现相比基于 Montgomery 的 NTT/INTT 系统的 C 实现高出 6.37~7.27 倍.

Pham 等人^[31]同样基于 Barret 算法提出了一种新的模乘架构,如图 4 所示,该架构由 2 个 DSP 计算单元和一个模块化约简过程(MR)组成.算法使用两个 DSP 对 23 位整数执行整数乘法,其中 MR 包括上半部分乘法器(UH-Mult)和下半部分乘法器(LH-Mult).UH-Mult 单元执行将 DSP 计算单元输出的高 24 位 V 与预计算常数 T 的相乘,而 LH-Mult 单元执行 UH-Mult 单元输出右移 24 位后与模数相乘,两个乘法过程被优化为移位、加法

和减法操作,优化后的 Barrett 模约简算法避免了使用乘法器进行大整数乘法,从而在硬件实现上显著降低了资源消耗和功耗.为了进一步提升硬件效率,Pham 等人^[31]在架构中部署了四级流水线的设计.第一和第二级位于 DSP 片内,在执行乘法和加法输出后,不消耗额外硬件资源,充分利用了 DSP 的内部结构,在不增加资源开销的情况下实现了流水线.第三级放置在 UH-Mult 的输出处,通过将长关键路径分割成若干个较短的路径,从而降低关键路径延迟,提高电路的最高工作频率.第四级位于模乘法器的输出处,进一步优化关键路径的时序表现,提高整个电路的工作频率. Pham 等人^[31]所设计的模乘法器架构能够以 385 MHz 的频率运行,整个模块仅消耗 111 个 LUTs、71 个 FFs 和 2 个 DSPs,与 Wang 等人^[32]的设计相比,硬件资源消耗减少到原来的 1/149,显著降低了硬件资源的使用.

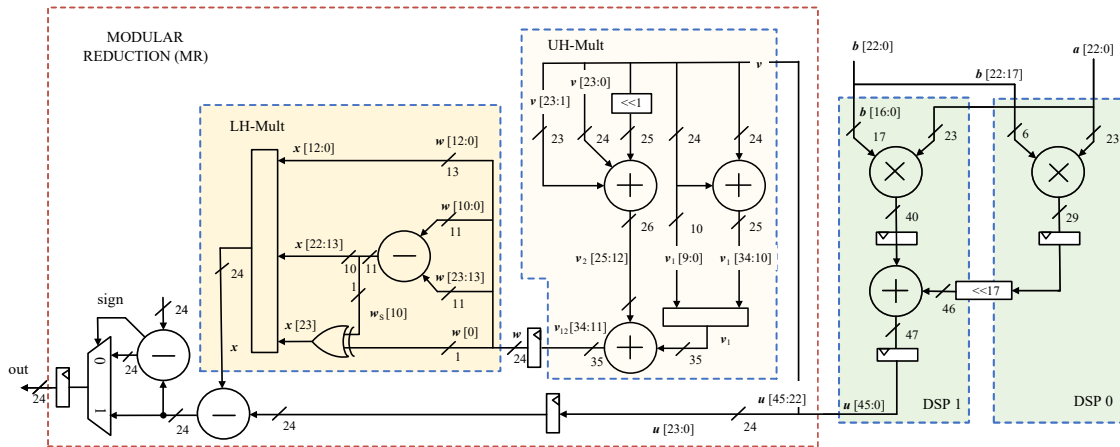


图 4 基于 Barret 约简算法的模乘法器架构^[31]

Malal 等人^[33]则选择了一种运行时可配置的 Montgomery 模约简算法,以增强设计的适应性和灵活性.该算法能够让设计的蝶形单元适应不同大小和类型的多项式操作,为 NTT 和 INTT 计算以及运行时的逐点乘法提供了灵活性,在资源使用和计算速度之间实现灵活的权衡.但是 Montgomery 模约简算法与其它算法相比,需要对被乘数进行预处理,以将它们转换到 Montgomery 域.这一预处理步骤需要额外的计算资源,会增加算法的总体资源消耗.

Li 等人^[25]的设计基于 Karatsuba 算法简化了模乘法算法,设计的模乘法单元可以计算 23 比特和 13 比特乘法.其设计基于 Karatsuba 算法可以减少关键路径的长度并提高系统的最大工作频率,进而提高了处理器的性能.然而 Karatsuba 算法通过将大整数乘法分解成较小整数乘法的方式,虽然减少了单个乘法操作的复杂度,但是增加了乘法器的个数,并且引入了额外的加法操作,这会增加系统的复杂性和延迟.此外,该设计存

在 DSP 资源浪费的情况,性能还可以进一步提升.在处理小规模数据时,Karatsuba 算法的优势不明显,甚至导致性能下降.因此,在实际应用中,需要综合考虑数据规模和系统需求,权衡 Karatsuba 算法的优点和缺点,以确定是否适合特定的硬件实现.

本节针对 CRYSTALS-Dilithium 算法中模乘操作,介绍了现有研究中的各种优化策略,这些策略涉及硬件结构、算法改进以及资源管理等多个方面.其中 Land 等人^[27]、Gupta 等人^[28]和 Zhao 等人^[22]的设计均是采用经典模约简策略递归利用模数性质进行约简, Land 等人^[27]使用大量 DSP 资源,实现了快速计算,并且通过移位寄存器延迟 DSP 输入以增加吞吐量,但存在 DSP 资源浪费的问题.

Gupta 等人^[28]在此基础上使用 LUTs 和 FFs 代替 DSPs,减少资源消耗. Zhao 等人^[22]进一步将模量转换为 CSD 表示法,压缩位宽,并通过多次的比较判别条件完成模约简,减少了数据位宽需求. Hwang 等人^[30]和

Pham 等人^[31]的设计基于 Barrett 算法. Hwang 等人^[30]通过结合整数近似的概念与 Barrett 算法,用于乘法计算之前对乘数进行约简来优化 NTT/INTT 的性能. Pham 等人^[31]基于 Barrett 算法提出了一种新的模乘架构可以在使用更少的硬件资源的同时实现相当的吞吐量. 该算法将模乘运算中的除法转化为乘法和移位操作,避免了直接执行昂贵的除法运算,从而降低了硬件复杂度和运算延迟,此外,Barrett 算法的各个步骤(如乘法、移位、减法)之间具有较好的并行性,可以通过流水线设计进一步提高运算效率和吞吐量. 但是 Barrett 算法在计算过程中需要预计算一个与模数相关的常数,这个预计算过程本身需要一定的计算量,当模数频繁变化时,预计算开销可能会抵消约简过程的性能优势,因此,对于某些特定的模数[如梅森(Mersenne)素数],基于 Barrett 算法的实现可能不如经典的 Montgomery 算法高效,例如 Malal 等人^[33]为了实现更好的灵活性和适应性,采用运行可配置的 Montgomery 算法. 但是考虑到 Montgomery 算法需要对乘数进行预处理,这在模数频繁变化时可能会成为性能瓶颈. 因此,研究者需要根据具体应用场景权衡不同算法的性能和实现成本来选择最合适的模约简算法.

3.2 NTT 的优化方案

NTT 在 Crystals-Dilithium 算法中扮演着核心角色,其性能直接影响到整个密码系统的效率. 对于 NTT 计算,一般采用 CT(Cooley-Tukey)蝶形单元应用在 NTT 过程中,而 GS(Gentleman-Sande)蝶形单元应用在 INTT 的过程中. 针对 NTT 的优化主要集中在硬件实现的结构优化和前文提到的模乘算法的改进上. 硬件实现的典型结构是使用蝶形单元按照蝶形图进行逐层计算^[34-36]. 目前,针对 NTT 的设计方案主要包括高速率设计^[22,27,37]和轻量级设计^[38].

针对 NTT 的高性能实现,基本采用并行设计结构,在同一层中同时计算多个蝶形单元^[39-41]. 如图 5 所示, Ricci 等人^[37]基于 FPGA 平台,提出了一种旨在充分利用硬件资源以实现高效 NTT 计算的设计方案. 他们设计了一个 2×2 NTT 蝶形计算单元的结构,并通过 4 个蝶形单元实现并行化和流水线化处理. 这样的设计合并了两个蝶形变换过程,有效减少了存储中间结果所需的块随机存取存储器(Block Random Access Memory, BRAM)数量,并减少了蝶形计算迭代间的间隔次数. 四个 BRAM 被用于存储输入值、中间结果以及输出值. 在未使用时,这些 BRAM 接口可以通过控制单元切换到 NTT 组件接口,以便进行数据的输入和输出. 在蝶形单元计算过程中, BRAM 接口被控制单元用于设置地址,同时连接到数据总线和 2×2 蝶形单元布局的计算结构中. 该设计方案同时考虑了并行计算和存储器访

问的优化,使其能够在不同硬件资源的 FPGA 平台上灵活部署,具有较强的可扩展性. 然而,这种 2×2 蝶形单元设计也存在局限性. 它限制了蝶形单元的复用,并且由于内存访问顺序的复杂性,需要一个复杂的控制模块来管理,这在 LUTs 和 DSPs 资源的消耗上有显著增加.

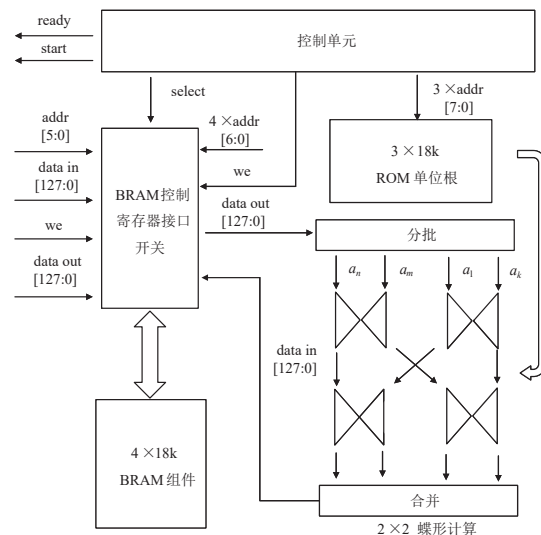


图5 2×2 NTT 蝶形计算单元^[37]

Land 等人^[27]在 Ricci 等人^[37]的 2×2 NTT 蝶形计算单元结构的基础上,选择使用 DSPs 来进行所有的算术操作,以实现低面积占用和高频率. 他们的设计特别利用了 FPGA 中 DSP 的不同操作模式,包括运行时重新配置、预添加以及单指令多数据(Single Instruction Multiple Data, SIMD)添加,这使得他们的 NTT 实现了 311 MHz 的最大频率,这一频率在当时基于 Xilinx Artix-7 平台的设计中是最高的,然而,这种高性能的实现也带来了非常高的 DSP 使用率. 此外, Land 等人^[27]充分发挥了 BRAM 模块的双端口功能,使他们的设计能够在 NTT 的最底层同时读取两个旋转因子,并同时处理四个系数的计算. 这种设计提高了数据吞吐量,但同时也意味着每个蝶形单元在每个周期中需要读取和回写两个系数,导致所需的 BRAM 数量与蝶形单元的数量成正比. 具体来说,如果有 k 个蝶形单元,就需要 k 倍的 BRAM 资源. 虽然 Land 等人^[27]的设计在频率上达到了显著优势,但也因为大量使用 DSP 资源和 BRAM 资源而造成了资源的利用率较低.

同样选择 2×2 NTT 蝶形计算单元进行 NTT 计算, Beckwith 等人^[38]提出了一种由地址解码器和 4 位先入先出队列(First In First Out, FIFO)组成的无冲突内存调用方式,如图 6 所示. 其中, 4x FIFO 组件用于转置 4×4 系数矩阵,以确保在 NTT 操作前和 INTT 操作后系数被正确排序. 地址解释器用于将 NTT 算法中的表示地址(RAddr)

出了一种 2×2 NTT蝶形计算单元结构,通过并行化和流水线化处理减少了BRAM的使用,并降低了计算迭代间的间隔次数.该设计优化了并行计算和存储器访问,提高了NTT计算的灵活性和可扩展性,但增加了LUTs和DSPs资源的消耗.Land等人^[27]利用FPGA DSPs资源进行所有算术操作,实现了低面积占用和高频率,通过运行时DSPs的重新配置、预添加和SIMD添加等操作模式,达到了311 MHz的高频率,但DSPs使用率极高,在资源消耗上较大.Beckwith等人^[38]提出了一种无冲突内存调用方式,通过地址解码器和FIFO组件优化了系数矩阵的转置和排序.该设计使用LUTs和FFs代替常数乘法,减少了DSPs的使用,但增加了LUTs和FFs资源的消耗.Zhao等人^[22]提出了分段流水式的NTT结构,用BRAM和计数器取代了大的移位寄存器,减少了内存访问和资源消耗.该设计通过优化内存访问顺序和减少BRAM消耗,提高了计算效率,但增加了NTT的关键路径.Land等人^[27]的设计虽然在频率上达到了优异的性能,但这种高性能的实现是以大量DSPs资源消耗为代价的.相比之下,Beckwith等人^[38]和Zhao等人^[22]通过优化内存管理和减少DSPs使用,实现了资源的更高效利用.可见,NTT的优化方案需要在性能提升、资源消耗和计算效率之间找到平衡,在追求高性能计算的同时如何平衡资源消耗将是未来研究开展的难点也是重点.

3.3 哈希和伪随机采样的优化方案

在CRYSTALS-Dilithium算法中,哈希函数和伪随机数生成器的作用是多方面的.采样过程与伪随机数生成紧密相关,它依赖于伪随机数生成器来从特定的分布中选择元素,这些元素用于构建多项式、矩阵等.优化这些操作可以显著提高算法的实现效率.

Crystals-Dilithium算法主要使用2015年被NIST标准化的第三代安全散列算法(Secure Hash Algorithm, SHA-3)(FIPS-202)^[44]家族中的SHAKE-128和SHAKE-256可扩展输出函数(XOFs).它们均是基于Keccak核设计的,统一使用1600-bit的状态(State)的Keccak_f[1 600]函数进行多轮f函数置换,每轮置换包括5种映射方式($\theta, \pi, \rho, \chi, \iota$),具体细节详见文献^[42].作为SHA-3函数的核心,Keccak核排列占据了超过20%的资源消耗和处理周期^[45],因此,Keccak核的优化设计是SHA-3函数

优化的关键.基于已有的研究成果,SHA-3函数的优化方案主要聚焦于Keccak核的轻量级或高速率设计.针对SHA-3函数的低资源优化设计,使用单个Keccak核实现所有功能,牺牲性能换取低资源消耗^[27,28].单个Keccak核的低资源优化方案主要采用折叠结构(Folded),其中折叠结构分为“通道式”(lane-wise)和“切片式”(slice-wise)两种折叠方式.Bertoni等人^[46]首次提出基于通道式的折叠方式,每次轮函数迭代处理1个lane(25bit),这种设计方式进行每轮5种函数映射时需要调用大量的前轮函数的数据,导致高延迟且每轮需要215个时钟周期数.Jungk等人^[47]提出了首个切片式的折叠方案,每次可以并行处理8个slice.切片式折叠允许在一个时钟周期内完成对折叠的完整轮次函数处理,但需要重新调度轮次函数.这种折叠方法是一种高度灵活的解决方案,关于并行处理的slice数量可以通过对控制逻辑进行小的修改而轻松调整.随后,Sundal等人^[48]和Jungk等人^[47]将Keccak_f[1 600]函数按照同样的slice-wise方式,一次并行处理8个或者16个slice,即每次迭代处理200 bit或者400 bit进行多轮置换,如图8和图9所示,相比Jungk等人^[47]的工作降低了延迟,达到了一定的硬件资源与计算效率的均衡.

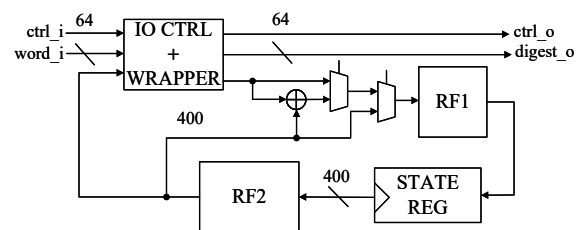


图8 每次处理400 bit的Keccak核折叠结构^[48]

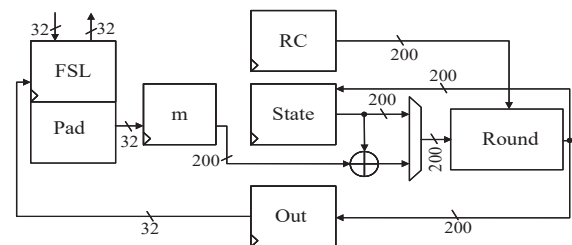


图9 每次处理200 bit的Keccak核折叠结构^[47]

表4 NTT的硬件实现对比

设计	最大工作频率/MHz	LUTs	FFs	DSPs	BRAMs	测试平台
文献[42](仅NTT乘法模运算单元)	152	7 073	2 031	0	NR	Artix-7
文献[43](仅NTT乘法模运算单元)	216	2 044	NR	16	NR	Artix-7
文献[33](资源节约型)	142.8	1 601	699	10	5	UltraScale+
文献[33](高吞吐量)	142.8	11 558	4 912	80	40	UltraScale+

针对 SHA-3 函数的高性能优化设计,则是通过级联多个 Keccak 核实现高速的计算,但是缺点是消耗大量的面积资源^[22,25,37,38]. 例如,Beckwith 等人^[38]在其设计方案中使用了三个 Keccak 核,其中两个主要用于 A 矩阵的采样,而第三个用于其余的哈希和采样. 为了提高采样性能,Beckwith 等人^[38]还实现了多个系数的并行采样,以充分利用 Keccak 核,但是这也会导致较大的资源消耗. 除此之外,也有少数研究在基于格的密码方案中对 SHA-3 函数进行了高速高效的硬件实现和优化:Wong 等人^[49]首次提出具有 2 个 IPR 和 2 个 OPR 两次展开架构,并且将 64 比特轮常数(Simple Round Constants, SRC)简化为 8 比特的方法. 在此基础上,刘冬生等人^[50]提出了一种可以适用于所有 SHA-3 函数的高效高速的硬件结构. 通过进一步压缩轮常数为 7 比特,并调整流水线寄存器的位置,进一步缩短关键路径,显著提升整个流水线架构的性能,提高了 33.4% 的吞吐量和 28.2% 的硬件效率,相关设计如图 10 所示.

但是这些采样算法都不支持在 32 位字节上工作. Land 等人^[27]设计了一个缓冲区来解决 32 位字与采样算法位宽不匹配的问题,实现了费舍尔-耶茨洗牌算法以优化挑战 c 的采样过程,并通过使用具有可变深度的移位寄存器来存储和处理非零系数的偏移量和符号位. 此外,作者还通过优化拒绝采样过程中的随机偏移量查找和替换操作,以及管理移位寄存器的深度来存储新的偏移量,提高了采样效率. 最终,处理后的多项式直接写入 BRAM,简化了整个数据处理流程,从而提升了算法的硬件实现效率和资源利用率. 虽然作者这样的处理可以解决哈希算法输出与采样算法位宽不匹

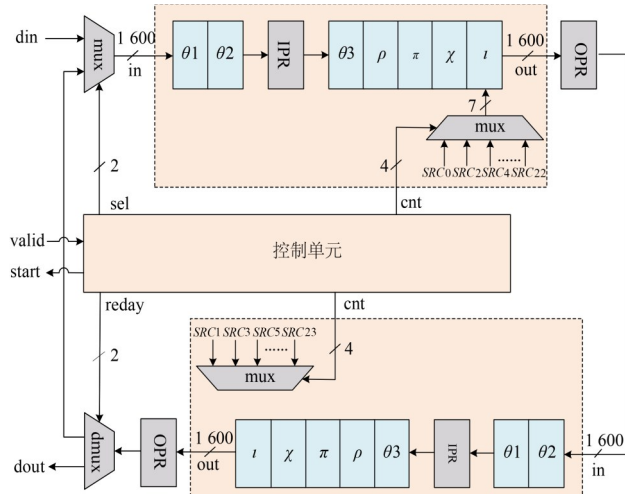


图 10 SHA-3 中 Keccak 核的高性能硬件结构^[50]

配的问题,但是会造成大量的寄存器的消耗,并且由于输出之后还需要格式转换,计算周期较长.

在此基础上,Zhao 等人^[22]提出了如图 11 所示的一种包含拒绝采样和基于双端口 BRAM 的 SampleInBall 采样模块,该模块使用 BRAM 阵列的空闲区域记录系数,代替了寄存器记录偏移量,从而避免了使用大量的寄存器. 其次,该模块的输出以标准多项式格式,可以直接进行 NTT 处理,无需额外的格式转换. 此外,直接的串行处理需要将每个有效样本进行一次读取操作和两次写入操作,对应于三个时钟周期,而提出的模块平均每个样本只需要一个周期. 其中黑色、蓝色和红色组件表示基于不同情况下的处理,实现存储和处理时间减少的同时,提高了系统的运行效率. 整体的设计达到低资源消耗的同时,减少了运算的周期.

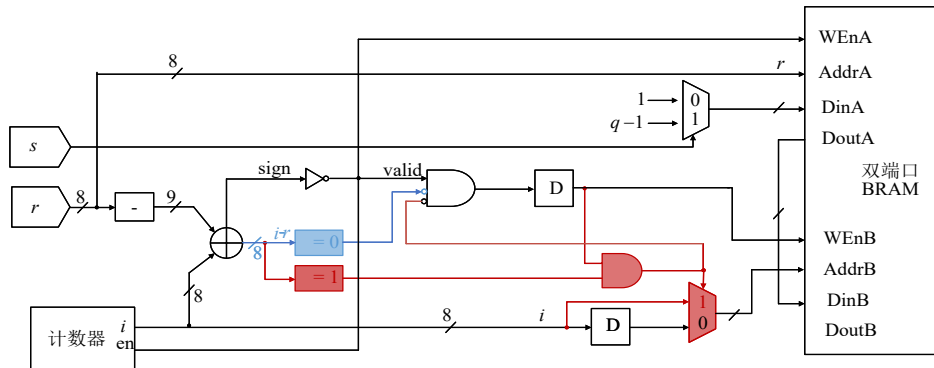


图 11 基于 BRAM 的采样单元^[22]

本节针对 CRYSTALS-Dilithium 算法中的哈希和采样操作的优化策略展开介绍,但是总的来说,目前对于哈希函数和采样还有巨大的优化空间. 尽管 SHA-3 函数的 FPGA 优化实现研究已经有了很多进展,但是只有少数基于 CRYSTALS-Dilithium 算法的文章使用相关的

优化策略. 目前根据 SHA-3 函数结构的优化方式可以分为:折叠结构、流水线结构和展开结构. 折叠结构虽然占用的资源小,但是运行周期长并且需要较为复杂的控制逻辑;展开结构是在一个时钟周期内实现 Keccak 核的多轮迭代,缺点是相应的资源消耗也会成倍增

了阶段间的依赖和数据传输,从而优化了性能和资源利用. 具体来说,签名过程被划分为预计算、阶段0和阶段1三个阶段. 预计算仅在签名生成开始时执行一次,涉及私钥的采样和解码,由第一组模块执行. 同时,另一组模块执行阶段0,包括 y 的采样和矩阵乘法,以生成签名尝试所需的 w 向量. 两个阶段完成后,最初用于预计算的模块组利用阶段0的结果尝试创建有效签名,即阶段1. 在执行阶段1的同时,另一组模块再次运行,以便在阶段1失败时立即进行新的尝试,形成了一个高效的两阶段流水线. 尽管这种设计提高了处理效率,但它仍然以牺牲资源消耗为代价.

参考以上的设计思路,Zhao 等人^[22]提出了一种分段流水线处理方法,在更短的操作周期和更少的存储需求下实现了高性能. 如图 14 所示,Zhao 等人^[22]将算法中的操作分为多个部分,不同的段被串行处理,段内的操作以流水线方式处理. 这种流水线方法减少了对中间结果的存储需求,并降低了内存访问次数;而分段处理则减少了算法所需的硬件资源,并允许更高效地利用 BRAM 模块为了进一步减少 BRAM 的存储消耗. 为了进一步降低 BRAM 的存储消耗,作者调整了预先计算和存储的矩阵 A , 改为即时计算模式. 此外,考虑到 NTT 计算与脉冲宽度调制(Pulse Width Modulation, PWM)点乘计算的紧密相关性,他们设计了 HEAD 模块和 TAIL 模块,分别放置在流水线中的 NTT 模块之前

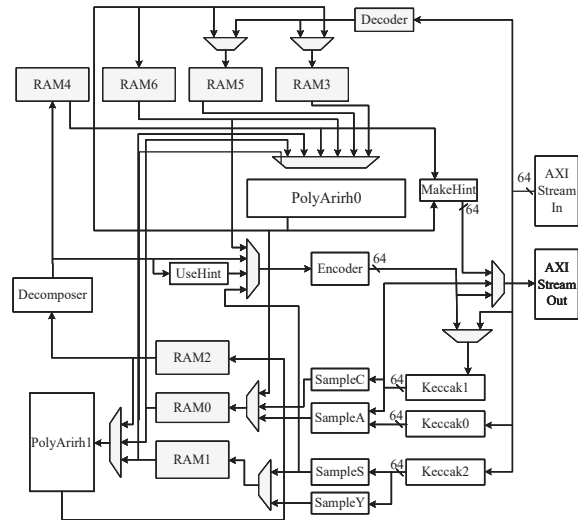


图 13 CRYSTALS-Dilithium 组合架构的块图^[38]

和之后,以实现高速流水线的多项式乘法. HEAD 模块集成了一个模乘法器和一个模加法器,而 TAIL 模块则包含比较器、模加法器、Decompose 模块等. 这种设计使得整个多项式乘法过程无须存储中间结果,从而实现了较高的处理速率. 然而,这种设计的缺点在于无法复用 NTT 资源,这在一定程度上导致了显著的面积开销. 这主要是因为多路复用的增加和额外模块实例的需求.

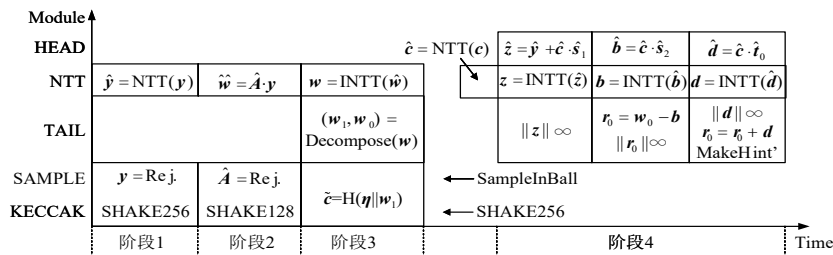


图 14 签名算法循环部分的分段流水线处理^[22]

本节介绍了 CRYSTALS-Dilithium 算法在追求高性能硬件实现时的不同策略和优化方案,这些策略涵盖了性能与资源消耗、并行化与流水线处理以及 BRAM 资源复用与优化. Ricci 等人^[37]的实现性能上取得了显著提升,但资源消耗较高. Beckwith 等人^[38]和 Zhao 等人^[22]通过优化设计减少了资源消耗,但这种优化在一定程度上牺牲了性能或增加了硬件复杂度. 其次,本节的研究都强调了并行化和流水线处理在提高性能中的重要性. Ricci 等人^[37]通过并行化处理提高了吞吐量,而 Beckwith 等人^[38]和 Zhao 等人^[22]通过并行化结合流水线处理减少了操作周期和存储需求. 此外, Zhao 等人^[22]的设计通过减少 BRAM 存储消耗和即时计算矩阵 A , 提高了处理速率,但增加了硬件的复杂度和面积开

销. 这种权衡在高性能硬件实现中也是常见的,尤其是在追求高吞吐量和快速响应的应用场景中.

4.2 低资源消耗的硬件实现架构

在面向资源受限的应用场景时,由于场景设备本身硬件资源十分有限,因此在设计的时候需要着重考虑低硬件资源消耗及低功耗进行轻量级的优化实现,如何在 LUTs、DSPs、FFs 三种不同的硬件资源中进行取舍从而获得最好的性能表现,并对核心模块进行复用,使得设计和实现的效率更高^[54],一直是一个研究难点.

Land 等人^[27]的设计侧重于在尽可能少的 LUTs 和 FFs 资源利用下实现可能的最佳性能. 作为第一个被提出用于低端 Xilinx Artix-7 平台的设计, Land 等人^[27]在他们的体系结构中通过有效使用 DSPs 以减少资源消

耗,复用各种模块,大量使用 DSPs 代替 LUTs 和 FFs 的面积,在一定程度上减少了 LUTs 和 FFs 资源的使用,但该体系结构的并行度较低,导致其模块的利用率较低,并且在设计过程中忽略了 DSPs 对于资源的占用. 其中该设计在采样模块使用单个共享实例来实现所有功能,也是为了节省资源而牺牲一些性能.

Gupta 等人^[28]的设计则选择使用 LUTs 和 FFs 代替 DSPs 从而减少资源消耗的优化方案,提出将协议算法中的 Power2Round 和 Decompose 两个分解运算模块如图 15 所示进行联合封装,共享两个模块的控制单元,并

且在 Decompose 模块的第二个乘法中预先计算了 16 个可能出现的值,使用一个 4 位 LUTs 作为输入优化计算资源,从而节省对于 DSPs 的使用. 其次,Gupta 等人^[28]还提出了基于 Barrett 算法的模约简算法以提高效率和减少资源使用,利用模约简操作有效地降低 DSPs 的数量. 此外,在采样模块,Gupta 等人^[28]选择节省资源而牺牲一些性能,使用单个共享实例来实现所有功能. 该设计在 UltraScale+ 平台上达到了 1.176 GHz 的运行频率,并且在整体设计仅需要使用 13 900 个 LUTs、6 800 个 FFs、4 个 DSPs 以及 35 个 BRAMs.

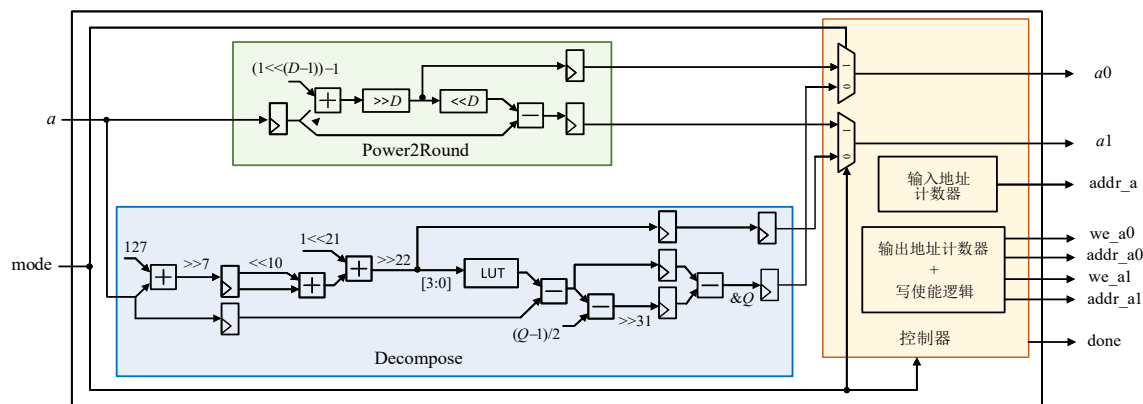


图 15 Power2Round 和 Decompose 联合封装模块^[28]

本节介绍了 CRYSTALS-Dilithium 算法在追求低资源消耗硬件实现时的不同策略和优化方案. 一个优化方案是用大量的 DSPs 单元进行大数乘法^[28],可以显著提高系统的频率,降低 LUTs 和 FFs 的资源消耗;另一个优化方法是通过更多地使用 LUTs 和 FFs 来替代 DSPs 的功能^[28],以降低对昂贵或能耗较高的 DSPs 资源的依赖. 选择哪种优化策略取决于具体的应用需求、成本预算、功耗限制以及性能要求.

4.3 后量子密码算法的联合实现

公钥加密、密钥协商(密钥封装)和数字签名技术作为公钥密码体制的三大重要组成部分,经常作为一个整体出现. 由于分组密码加密速度比公钥加密速度快,通常在密码硬件设备中均会集成分组密码算法、哈希算法、密钥封装算法和数字签名算法,这些算法足以满足多种应用场景下的安全防护要求. 在硬件实现过程中,一些密码算法会共用某些电路组件,如密码封装算法 Crystals-Kyber 和数字签名算法 Crystals-Dilithium 共用 NTT 模块等. 所以,研究者们通过在密码算法实现中有效地复用的通用核心算子电路,提升整体资源利用率. 具体的联合实现优化方法包含:Crystals-Kyber 和 Crystals-Dilithium 的联合高性能硬件实现以及 Saber 和 Crystals-Dilithium 的联合高性能硬件实现.

Crystals-Dilithium 和 Crystals-Kyber 均利用 NTT 用

于高效地执行多项式乘法,二者的联合实现可以满足高安全性并适配多种场景. 运用 Crystals-Kyber 对发送的信息进行公钥加密,同时运用 Crystals-Dilithium 对发送和接收端身份信息进行数字签名加密. 现阶段研究的难点主要有两点:(1)Crystals-Kyber 中使用的是 13 比特多项式模乘,而 Crystals-Dilithium 使用的是 23 比特多项式模乘. 二者不同位数的乘法使得在联合实现中,多项式模乘法单元的设计非常困难. 同时,二者不同的模数也会带来不同的模约简要求.(2)Crystals-Kyber 和 Crystals-Dilithium 整体实现流程不同,要求在设计的时候选择合适的内存调用方式、逻辑控制单元以及尽可能寻求更多的模块复用,才能达到性能和资源的平衡.

Aikata 等人^[55]首次针对 Crystals-Dilithium 和 Crystals-Kyber 提出了一种统一的体系结构 KaLi,该体系设计了一个统一而灵活的多项式算术单元,可以为 Crystals-Dilithium 和 Crystals-Kyber 的所有安全级别执行密钥生成、封装、解封装、签名生成和签名验证,并且其处理 Crystals-Kyber 运算的速度是 Crystals-Dilithium 运算的两倍. Aikata 等人^[55]提出了一种用于 Kali 的算法,使为 Crystals-Dilithium 设计的整数乘法单元对 Crystals-Kyber 具有灵活性,它执行两次 23 比特乘 12 比特的整数乘法,如图 16 所示. 这个算法在 Crystals-

Dilithium 的情况下给出了一个乘系数,在 Crystals-Kyber 的情况下给出了两个乘系数,如图 17 所示。

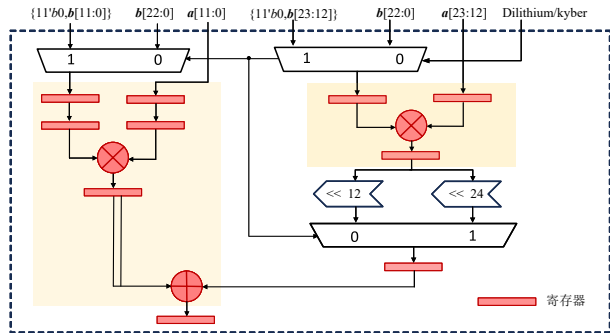


图 16 适用于多模数的整数乘法器^[55]

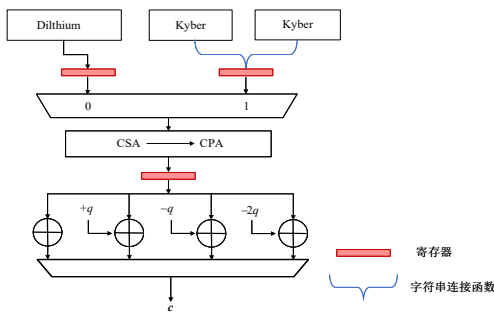


图 17 Crystals-Kyber 和 Crystals-Dilithium 统一的模约简单元^[55]

KaLi 设计在使用多个时钟域的 ASIC 平台上特别适用。在 ASIC 的 28 nm/65 nm 技术上,该设计实现了 2 GHz/560 MHz 的时钟频率,可以用于存储单元的快速时钟。在 Zynq Ultrascale+ 平台上,其所提出的架构可以在 270 MHz 时钟频率下仅仅使用 23 277 个 LUTs、9 758 个 FFs、4 个 DSPs 和 24 个 BRAMs。经验证,KaLi 方法所达到的性能比这两种方案中的任何一种的独立实现都要好。

Li 等人^[42]提出一种 2KNTT 多项式乘法结构,如图 18 所示,实现模乘法单元可以同时满足不同模数的操作,并且控制单元可以满足不同项数的要求,同时适用于 Crystals-Dilithium 和 Crystals-Kyber 算法。在原始的 NTT 算法中,模数 q 需要满足 $q=1 \bmod 2n$ 的条件,而在 2KNTT 算法中,这种限制被削弱为只需要满足 $q=1 \bmod n$ 或 $q=1 \bmod n/2$,甚至更小的条件。这种改进的算法使得在加密方案中选择更小的模数成为可能,从而减少了带宽需求和资源消耗。在 Artix-7 平台以到 152 MHz 的频率运行,消耗 7 073 个 LUTs、2 310 个 FFs 和 0 个 DSPs。

Li 等人^[42]改进的 2KNTT 算法的多项式乘法单元结构如图 18 所示,由于所选算法本身的时间复杂度较低,计算和存储体系结构采用流水线型的设计,实现了运算时间和面积的平衡,而整体的并行结构和统一的蝶形计算单元则保证了性能的提高。与高度并行设计相

比,它平衡了高并行条件下计算单元统一和存储单元调度的复杂问题,而更高的性能可以降低单位性能所需的资源量。

虽然 Saber 在第三轮中因为特定情况下的安全性难以被证明而被淘汰出局,但是 Saber 相对于其他格密码方案的实现较为容易,主要的原因是 M-LWR 所决定的采样方式和独特的基于 2 的幂次的模数,使得在设计硬件实现方案时结构较为简单,减少了随机性以及所需带宽。因此 Saber 在联合实现方面仍然具有较高的研究价值。Basso 等人^[43]和 Aikata 等人^[56]提出了可以同时实现于 Crystals-Dilithium 数字签名以及 Saber 密钥封装机制的联合设计。

除了 Crystals-Kyber 和 Saber 与 Crystals-Dilithium 的联合实现外,椭圆曲线数字签名算法(Elliptic Curve Digital Signature Algorithm, ECDSA)作为一种基于椭圆曲线密码学的数字签名算法,也有与 Crystals-Dilithium 混合的数字签名方案, Ghinea 等人^[57]提出了一种 ECDSA 和 Crystals-Dilithium 混合的数字签名方案软件实现,以在经典加密算法和后量子安全方案之间取得平衡。这种方法为当前的安全密钥提供了一种过渡方案,既保持了经典加密算法的安全性,又具有抵抗量子攻击的能力,使其能够在量子计算时代仍然保持安全。Ghinea 等人^[57]采用了一种称为“Strong Nesting”的方法,首先使用 X-EUF-CMA 安全方案对给定的消息进行签名,得到 σ_1 ,然后使用 X-SUF-CMA 安全方案对 (m, σ_1) 进行签名,得到 σ_2 ,最终的混合签名为 $\sigma=(\sigma_1, \sigma_2)$ 。该方法保持了 Crystals-Dilithium 的强不可伪造性,同时也保持了 ECDSA 的存在性不可伪造性,即使一个方案被攻破,也可以保持另一个底层方案的安全性,在安全密钥身份验证的情境下具有强不可伪造性,因此对经典和量子攻击都具有抵抗力。可惜该方案目前还停留在软件实现层面,在未来可以在硬件实现层面做深入研究。

5 Crystals-Dilithium 硬件实现对比分析

本节将结合第 3 节和第 4 节的关键组件优化与整体架构优化方案,全面对比基于这些优化方案的 Crystals-Dilithium 硬件实现的性能。通过文字描述和表格的形式,我们将展示不同优化策略对性能的影响,以便更好地理解各方案在实际应用中的效果。

表 5 汇总了已有的 Crystals-Dilithium 协议的硬件实现,从表中可以看出,大部分设计使用的 DSP 数量是非常少的,这是迭代实现的一个优点。迭代实现的另外一个优点是能够节省很多 LUTs 和 FFs 的资源消耗。由于迭代型模多项式乘法模块只有一个蝶形计算单元,所以需要大量的 BRAM 资源,一方面暂存中间计算数据,另一方面需要存储逻辑地址和计算地址来控制 NTT 计

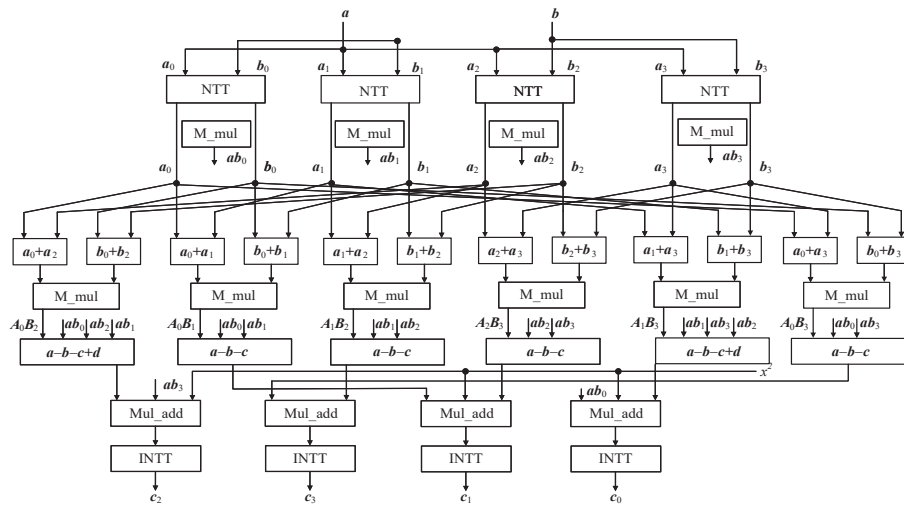


图 18 2 KNTT 多项式乘法结构^[42]

算的顺序. Beckwith 等人^[38]、Land 等人^[27]、Aikata 等人^[55,56]的设计方案 BRAM 数量都超过了 20,其原因就是在于二者都使用了大量的 BRAM 进行 NTT 计算顺序的逻辑地址的控制. 其流水线型模多项式乘法模块虽然会消耗大量的 DSP 资源以及 LUTs、FFs 资源,但是在外围计算电路能够匹配乘法单元频率的前提下,流水线型设计仍能获得最好的运算性能表现,尤其是频率表现. 同时,部分设计方案采用流水线的设计方法,数据不需要在计算过程中间进行任何暂存和滞留,也不需要逻辑地址控制 NTT 计算顺序,所以 BRAM 的使

用量几乎为 0.

在高速率的硬件实现方案中, Ricci 等人^[37]设计了四个 2×2 的 NTT 蝶形计算单元实现并行化和流水线化的处理,使用消耗大量资源直接将算法映射到硬件上,虽然提高了硬件的性能,但是硬件资源的使用效率较低并且硬件消耗资源大. 在此基础上, Beckwith 等人^[38]提出一个能够在所有安全级别执行三种操作的组合模块,也包括每个模块支持一种操作的单个模块. 其优点在于通过组合模块的调用,三种操作可以在同一硬件电路中实现,节约了近 52% 的硬件资源开销,并且通过

表 5 硬件实现方案性能参数比较(安全级别 2/3/5)

设计	CCs/k	Time/ μ s	Freq/MHz	LUT	FF	BRAM	DSP	ATP	
文献[25] (高速率)	4.4/30.1/5.0	33.9/231/38.5	130 (Artix-7)	32 320	9 734	14	24	4.01	
	7.3/49.5/8.0	56.1/381/61.6							6.59
	10.9/55.2/12.7	83.8/425/97.6							8.02
文献[27] (高速率)	NR	18.7/9.6/66.9	163 (Artix-7)	51 028	27 042	53	45	NR	
		33.1/18.1/10.5							
		51.0/33.8/112.1							
文献[37] (高速率)	12.6/18.4/10.5	27.8/17.7/15	350 (UltraScale+)	184 372	146 494	178	316	NR	
	18.2/21.0/15.3	19.2/15.5/10.5							
	22.9/22.3/20.2	15.2/14.2/7.8							
文献[38] (高速率)	4.8/6.5/10.9	42/57/94	116 (Artix-7)	53 187	28 318	29	16	NR	
	8.2/9.7/16.7	71/84/139							
	14.1/14.6/24.4	121/126/210							
文献[22] (低资源)	4.2/28.1/4.4	43.3/290/45.4	96.9 (Artix-7)	29 998	10 366	11	10	4.04	
	5.9/44.7/6.1	60.9/461/63.0						6.23	
	8.8/49.0/9.0	90.1/506/92.9						7.34	
文献[28] (低资源)	14.6/31.7/15.4	54.1/117/57.0	270 (UltraScale+)	23 347	9 798	24	4	1.83	
	23.6/48.7/26.1	87.4/180/96.7						2.92	
	39.7/70.2/46.7	147/260/173						4.65	

FIFO 的使用对系统内存模块做了优化,实现更低的延迟时间和更小的面积资源.但因为该设计要兼容 FALCON 算法的电路需求,所以 Crystals-Dilithium 的性能受到了限制.参考这样的设计思路,Zhao 等人^[22]也提出了一个支持三个操作(密钥生成、签名和验签)以及所有安全级别的紧凑高性能硬件架构.其优点在于:

(1) 提出分段流水式的处理,减少时钟周期,实现数据独立执行的并行操作,这样的优化方式减少了大量操作的执行时间和中间结果的存储需求,解决了硬件资源利用率低的问题.

(2) 设计了包括高速流水线 NTT 模块,基于 BRAM 的 SampleInBall 采样模块,紧凑的分解模块和三个优化的模约简模块,仅需要使用 11 个 BRAMs 单元,减少了硬件资源开销.但是多路复用的增加,额外的模块实例等,会导致一定的面积资源消耗,同时受限于分段流水的结构以及架构中对 DSP 单元的利用率不高的原因,其最大工作频率为 96.6 MHz,频率性能还没有达到最佳.

在低资源消耗的硬件实现方案中,Land 等人^[27]的整体实现方案中,该方案使用了 45 个 DSPs 使得整体频率能够达到 205 MHz,但由于其结构单一,缺少对 DSP 单元的优化设计,造成了 DSP 资源的浪费,理论上该结构可以通过优化 DSP 单元的结构来获得更好的性能表现.Gupta 等人^[28]利用资源和控制逻辑共享、预计算的模块共享以及模块共享等多种优化策略,降低 LUTs 和 FFs 使用量的同时节省了对 DSP 单元的消耗.与 Beckwith 等人^[38]的设计相比,LUTs 的资源消耗减少了 73.7%,FFs 的资源消耗减少了 75.8%,DSPs 的资源消耗减少了 75%.

在联合实现的方案中,Li 等人^[42]和 Aikata 等人^[55]提出了 Crystals-Dilithium 和 Crystals-Kyber 的联合实现.Li 等人^[42]提出了 2KNTT 算法的多项式乘法架构,该算法时间复杂度较低,计算和存储体系结构的流线型控制复杂度决定了时间和面积的平衡,整体并行结构和统一的蝶形计算单元大大提高电路性能,在 Artix-7 平台上,NTT 模块的最大工作频率为 152 MHz,电路消耗的资源为 7 073 个 LUTs、2 310 个 FFs. Aikata 等人^[55]则提出了全新的 Kali 设计,设计了一种统一而灵活的多项式运算单元,联合实现了 Crystals-Kyber 和 Crystals-Dilithium 的所有计算过程,电路在 DSPs 的资源消耗上大大降低的同时,在性能表现上也较为优秀,电路消耗的资源为 23 000 个 LUTs、9 700 个 FFs、4 个 DSPs 和 24 个 BRAMs.

Basso 等人^[43]和 Mert 等人^[56]提出了 Crystals-Dilithium 和 Saber 的联合实现.由于 Saber 和 Crystals-Dilithium 共享相同的模数,使得模乘法单元的设计变得简单.但受限

于 Saber 的性能,使得二者的联合设计只能用于特定的应用场景. Basso 等人^[43]的设计在 Artix-7 平台上实现了使用了 333 MHz 的运行频率,消耗 2 012 个 LUTs、331 个 FFs 以及 6 个 DSPs. 而 Aikata 等人^[56]的设计在 Ultra-scale+ 平台上实现了 200 MHz 的运行频率,消耗 18 406 个 LUTs、9 323 个 FFs、4 个 DSPs 以及 24 个 BRAMs.

6 结束语

本文围绕 PQC 数字签名算法 Crystals-Dilithium 的国内外研究成果进行综述.针对 Crystals-Dilithium 数字签名可以根据不同的设计要求调整速率和资源之间的平衡的特点,首先,本文就模乘、NTT、哈希和采样等底层关键技术的优化方案进行了分析,介绍了底层关键组件的优化方向和优化策略;随后,针对目前已有的 Crystals-Dilithium 算法的整体实现,将其分为高速率、低资源消耗和算法联合实现这三个场景,对不同场景下的硬件实现的优化方法和硬件实现的结果进行了归纳总结和对比分析,展现了不同设计思路和硬件架构的优劣.

当前 PQC 算法标准化正如火如荼地展开,未来 PQC 迁移将是全球安全战略中的重要一环. PQC 算法的硬件实现相比软件实现具有更高的性能、更低的延迟、更高的并行性、更低的功耗、更高的安全性等优势,能够在各行各业中提供更高效、更安全的解决方案.所以,在量子计算迅速发展的时代,PQC 算法的硬件实现与优化设计研究愈发重要.强烈建议未来对于后量子密码算法的硬件实现研究可以集中在以下四个方面:

(1) 高速率硬件实现方案.通过改进多项式模乘算法和模约简算法减少乘法计算的时钟周期,并且通过改进硬件实现方案,改善并行度和内存调用逻辑等方式,提高整体硬件实现速率.

(2) 轻量级硬件实现方案,通过复用 NTT、keccak 核等多个核心组件,采用流水线的方式,减少存储需求和处理时间,降低硬件实现的资源消耗.

(3) 多种后量子密码方案联合硬件实现,当然可以不局限于 PQC 算法之间,也可以将 PQC 算法与经典密码算法进行联合实现.

(4) 高性能硬件实现方案,在速率和资源消耗之间寻找一种平衡,通过观察某些算法之间的组件关联度,对于通用的模块进行复用,展开电路设计与优化,减少不必要的硬件开销,为多应用场景提供全面的安全保障.最后,希望本综述可以为研究者在优化 PQC 算法的硬件实现方面提供思路和参考价值.

参考文献

[1] RIVEST R L, SHAMIR A, ADLEMAN L. A method for ob-

- taining digital signatures and public-key cryptosystems[J]. *Communications of the ACM*, 1978, 21(2): 120-126.
- [2] DAN B. The decision Diffie-Hellman problem[M]//*Algorithmic Number Theory*. Berlin, Heidelberg: Springer, 1998: 48-63.
- [3] ELGAMAL T. A public key cryptosystem and a signature scheme based on discrete logarithms[J]. *IEEE Transactions on Information Theory*, 1985, 31(4): 469-472.
- [4] KOBLITZ N. Elliptic curve cryptosystems[J]. *Mathematics of Computation*, 1987, 48(177): 203-209.
- [5] GALBRAITH S D. *Mathematics of Public Key Cryptography*[M]. Cambridge: Cambridge University Press, 2012.
- [6] SHOR P W. Algorithms for quantum computation: Discrete logarithms and factoring[C]//*Proceedings 35th Annual Symposium on Foundations of Computer Science*. Piscataway: IEEE, 1994: 124-134.
- [7] CHEN L, JORDAN S, et al. Report on Post-Quantum Cryptography[M]. Gaithersburg: US Department of Commerce, National Institute of Standards and Technology, 2016.
- [8] MICCIANCIO D, REGEV O. *Lattice-Based Cryptography*[M]//*Post-Quantum Cryptography*. Berlin, Heidelberg: Springer, 2009: 147-191.
- [9] BUCHMANN J, DAHMEN E, SZYDLO M. Hash-based digital signature schemes[M]//*Post-Quantum Cryptography*. Berlin, Heidelberg: Springer, 2009: 35-93.
- [10] OVERBECK R, SENDRIER N. *Code-based cryptography*[M]//*Post-Quantum Cryptography*. Berlin, Heidelberg: Springer, 2009: 95-145.
- [11] PETZOLDT A, BULYGIN S, BUCHMANN J. A multivariate based threshold ring signature scheme[J]. *Applicable Algebra in Engineering, Communication and Computing*, 2013, 24(3): 255-275.
- [12] FEO L D. Mathematics of isogeny based cryptography[EB/OL]. (2017-11-11)[2025-03-19]. <https://arxiv.org/abs/1711.04062>.
- [13] BEULLENS W. Breaking rainbow takes a weekend on a laptop[M]//*Advances in Cryptology-CRYPTO 2022*. Cham: Springer Nature Switzerland, 2022: 464-479.
- [14] MAILLARD B. Implementing the Dual Approaches for Solving LWE[D]. Lund: Lund University, 2023.
- [15] MAINO L, MARTINDALE C, PANNY L, et al. A Direct key recovery attack on SIDH[M]//*Advances in Cryptology-EUROCRYPT 2023*. Cham: Springer Nature Switzerland, 2023: 448-471.
- [16] DUCAS L, KILTZ E, LEPOINT T, et al. CRYSTALS-dilithium: A lattice-based digital signature scheme[J]. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2018, 2018: 238-268.
- [17] LANGLOIS A, STEHLÉ D. Worst-case to average-case reductions for module lattices[J]. *Designs, Codes and Cryptography*, 2015, 75(3): 565-599.
- [18] LIU W Q, FAN S L, KHALID A, et al. Optimized schoolbook polynomial multiplication for compact lattice-based cryptography on FPGA[J]. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2019, 27(10): 2459-2463.
- [19] SINHA ROY S, BASSO A. High-speed instruction-set coprocessor for lattice-based key encapsulation mechanism: Saber in hardware[J]. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2020: 443-466.
- [20] BERNSTEIN D J. Fast multiplication and its applications[J]. *Algorithmic Number Theory*, 2008, 44(1): 325-384.
- [21] WEIMERSKIRCH A, PAAR C. Generalizations of the Karatsuba algorithm for efficient implementations [EB/OL]. (2006-07-03)[2025-03-19]. <https://eprint.iacr.org/2006/224>.
- [22] ZHAO C K, ZHANG N, WANG H N, et al. A compact and high-performance hardware architecture for CRYSTALS-dilithium[J]. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2021, 2021: 270-295.
- [23] 芮康康, 王成华, 范赛龙, 等. 一种高性能R-LWE格加密算法的电路结构及其FPGA实现[J]. *数据采集与处理*, 2019, 34(4): 689-696.
- RUI K K, WANG C H, FAN S L, et al. High performance hardware architecture of lattice-based cryptography and its FPGA implementation[J]. *Journal of Data Acquisition and Processing*, 2019, 34(4): 689-696. (in Chinese)
- [24] MCIVOR C, MCLOONE M, MCCANNY J V. Modified Montgomery modular multiplication and RSA exponentiation techniques[J]. *IEE Proceedings - Computers and Digital Techniques*, 2004, 151(6): 402-408.
- [25] LI X, LU J H, LIU D S, et al. A high speed post-quantum crypto-processor for crystals-dilithium[J]. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 2024, 71(1): 435-439.
- [26] PLANTARD T. Efficient word size modular arithmetic[J]. *IEEE Transactions on Emerging Topics in Computing*, 2021, 9(3): 1506-1518.
- [27] LAND G, SASDRICH P, GÜNEYSU T. A hard crystal-implementing dilithium on reconfigurable hardware[M]//*Smart*

- Card Research and Advanced Applications. Cham: Springer International Publishing, 2022: 210-230.
- [28] GUPTA N, JATI A, CHATTOPADHYAY A, et al. Lightweight hardware accelerator for post-quantum digital signature CRYSTALS-dilithium[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2023, 70(8): 3234-3243.
- [29] HARTLEY R I. Subexpression sharing in filters using canonic signed digit multipliers[J]. IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing, 1996, 43(10): 677-688.
- [30] HWANG V, KIM Y, SEO S C. Barrett multiplication for dilithium on embedded devices[EB/OL]. (2023-12-25)[2025-03-19]. <https://eprint.iacr.org/2023/1955>.
- [31] PHAM T X, DUONG-NGOC P, LEE H. An efficient unified polynomial arithmetic unit for CRYSTALS-dilithium[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2023, 70(12): 4854-4864.
- [32] WANG T F, ZHANG C, CAO P, et al. Efficient implementation of dilithium signature scheme on FPGA SoC platform[J]. IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 2022, 30(9): 1158-1171.
- [33] MALAL A. Designing efficient and flexible NTT accelerators [EB/OL]. (2024-09-15)[2025-03-19]. <https://eprint.iacr.org/2023/1617>.
- [34] FRITZMANN T, SEPULVEDA J. Efficient and flexible low-power NTT for lattice-based cryptography[C]//2019 IEEE International Symposium on Hardware Oriented Security and Trust. Piscataway: IEEE, 2019: 141-150.
- [35] JATI A, GUPTA N, CHATTOPADHYAY A, et al. SPQ-Cop: Side-channel protected post-quantum cryptoprocessor[C]//IACR Cryptology ePrint Archive. Trier: DBLP, 2019: 765.
- [36] FRITZMANN T, SHARIF U, MÜLLER-GRITSCHNEDER D, et al. Towards reliable and secure post-quantum coprocessors based on RISC-V[C]//2019 Design, Automation & Test in Europe Conference & Exhibition. Piscataway: IEEE, 2019: 1148-1153.
- [37] RICCI S, MALINA L, JEDLICKA P, et al. Implementing CRYSTALS-dilithium signature scheme on FPGAs[C]//Proceedings of the 16th International Conference on Availability, Reliability and Security. New York: ACM, 2021: 1-11.
- [38] BECKWITH L, NGUYEN D T, GAJ K. High-performance hardware implementation of lattice-based digital signatures[EB/OL]. (2022-02-25)[2025-03-19]. <https://eprint.iacr.org/2022/217>.
- [39] MERT A C, OZTURK E, SAVAS E. Design and implementation of a fast and scalable NTT-based polynomial multiplier architecture[C]//2019 22nd Euromicro Conference on Digital System Design. Piscataway: IEEE, 2019: 253-260.
- [40] ZHANG N, YANG B H, CHEN C, et al. Highly efficient architecture of NewHope-NIST on FPGA using low-complexity NTT/INTT[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2020, 2020(2): 49-72.
- [41] FRITZMANN T, SIGL G, SEPÚLVEDA J. RISQ-V: Tightly coupled RISC-V accelerators for post-quantum cryptography[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2020, 2020(4): 239-280.
- [42] LI H Q, CHEN T, WU A Q, et al. High efficient architecture of polynomial multiplier with variable parameter based on 2KNNTT[C]//2022 IEEE Asia Pacific Conference on Circuits and Systems. Piscataway: IEEE, 2022: 606-610.
- [43] BASSO A, AYDIN F, DINU D, et al. Where star wars meets star trek: Saber and dilithium on the same polynomial multiplier[EB/OL]. (2022-03-08)[2025-03-19]. <https://eprint.iacr.org/2021/1697>.
- [44] NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY. SHA-3 Standard: Permutation-based hash and extendable-output functions[S/OL]. Gaithersburg: NIST, 2015(2015-08-04)[2025-03-19]. https://www.nist.gov/publications/sha-3-standard-permutation-based-hash-and-extendable-output-functions?pub_id=919061.
- [45] CARRIL X, KARDARIS C, RIBES-GONZÁLEZ J, et al. Hardware acceleration for high-volume operations of CRYSTALS-kyber and CRYSTALS-dilithium[J]. ACM Transactions on Reconfigurable Technology and Systems, 2024, 17(3): 1-26.
- [46] BERTONI G, DAEMEN J, PEETERS M, et al. Keccak[C]//Annual International Conference on the Theory and Applications of Cryptographic Techniques. Berlin, Heidelberg: Springer, 2013: 313-314.
- [47] JUNGK B, APFELBECK J. Area-efficient FPGA implementations of the SHA-3 finalists[C]//2011 International Conference on Reconfigurable Computing and FPGAs. Piscataway: IEEE, 2011: 235-241.
- [48] SUNDAL M, CHAVES R. Efficient FPGA implementation of the SHA-3 hash function[C]//2017 IEEE Computer Society Annual Symposium on VLSI. Piscataway:

- IEEE, 2017: 86-91.
- [49] WONG M M, HAJ-YAHYA J, SAU S, et al. A new high throughput and area efficient SHA-3 implementation[C]// 2018 IEEE International Symposium on Circuits and Systems. Piscataway: IEEE, 2018: 1-5.
- [50] 刘冬生, 陈勇, 熊思琦, 等. 应用于后量子密码的高速高效 SHA-3 硬件单元设计[J]. 信息安全学报, 2021, 6(6): 32-39.
- LIU D S, CHEN Y, XIONG S Q, et al. Design of high-speed and high-efficiency SHA-3 hardware unit for post-quantum cryptography[J]. Journal of Cyber Security, 2021, 6(6): 32-39. (in Chinese)
- [51] CONTI F, SCHILLING R, SCHIAVONE P D, et al. An IoT endpoint system-on-chip for secure and energy-efficient near-sensor analytics[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2017, 64(9): 2481-2494.
- [52] BANERJEE U, PATHAK A, CHANDRAKASAN A P. 2.3 an energy-efficient configurable lattice cryptography processor for the quantum-secure Internet of Things[C]// 2019 IEEE International Solid-State Circuits Conference. Piscataway: IEEE, 2019: 46-48.
- [53] DUCAS L, LEPOINT T, LYUBASHEVSKY V, et al. Cryptals-dilithium: Digital signatures from module lattices[EB/OL]. (2018-12-01) [2023-12-10]. <https://eprint.iacr.org/2017/633>.
- [54] GRECONICI D O C, KANNWISCHER M J, SPRENKELS D. Compact dilithium implementations on cortex-M3 and cortex-M4[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2021, 2021(1): 1-24.
- [55] AIKATA A, MERT A C, IMRAN M, et al. KaLi: A crystal for post-quantum security using kyber and dilithium[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2023, 70(2): 747-758.
- [56] AIKATA A, MERT A C, JACQUEMIN D, et al. A unified cryptoprocessor for lattice-based signature and key-exchange[J]. IEEE Transactions on Computers, 2022, 72(6): 1568-1580.
- [57] GHINEA D, KACZMARCZYCK F, PULLMAN J, et al. Hybrid post-quantum signatures in hardware security keys[M]//Applied Cryptography and Network Security Workshops. Cham: Springer Nature Switzerland, 2023: 480-499.

作者简介



崔益军 男, 1988年8月出生于江苏省海安市. 现为南京航空航天大学集成电路学院副院长、副教授、博士生导师. 研究方向为集成电路设计, 包括基于物理不可克隆函数的安全认证技术、面向未来信息系统的后量子密码安全通信技术、硬件木马、新兴计算等. 中国电子学会会员编号: E190036763M. E-mail: yijun.cui@nuaa.edu.cn



李梦雪 女, 2001年11月出生于浙江省温州市. 现为南京航空航天大学集成电路学院硕士研究生. 研究方向为后量子密码算法中核心算子的硬件设计与优化. E-mail: limengxue@nuaa.edu.cn



王 斐 女, 1991年2月出生于安徽省宿州市. 现为南京航空航天大学集成电路学院讲师. 研究方向为量子计算对公钥密码算法的安全性分析和后量子密码算法的软硬件加速研究. 发表SCI和EI论文10余篇. E-mail: wangbei91@nuaa.edu.cn



王成华 男, 1963年10月出生于江苏省扬中市. 现为南京航空航天大学集成电路学院教授、博士生导师. 研究方向为集成电路设计. E-mail: chwang@nuaa.edu.cn



刘伟强 男, 1983年3月出生于山东省东营市. 现为南京航空航天大学集成电路学院执行院长、教授、国家自然科学基金杰出青年基金项目入选者. 研究方向为高能效高安全新兴计算芯片. 中国电子学会会员编号: E190011136S. E-mail: chwang@nuaa.edu.cn