

Crystal-Kyber算法的FPGA高效并行优化

吕顺森^{1,2}, 李 斌^{1*}, 翟嘉琪^{1,2}, 李松岐^{1,2}, 周清雷¹

(1. 郑州大学计算机与人工智能学院, 河南郑州 450001; 2. 数字工程与先进计算国家重点实验室, 河南郑州 450001)

摘要: 多项式乘法运算制约着基于格的后量子密码在现实中的应用. 为提高后量子密码Crystal_Kyber算法的性能效率, 减少运行时间, 降低多项式乘法的影响, 本文设计了一种新的蝶形运算单元对素模 $q=3\ 329$ 的Kyber方案进行优化. 首先, 采用16路并行调度新型蝶形运算单元的方式执行算法, 缩短了计算周期; 其次, 使用流水线技术以及改进的K²RED算法, 设计实现新型蝶形运算单元, 用于降低资源消耗; 最后, 利用多RAM的方式存储数据, 并且多通道优化RAM, 允许数据交替存储在RAM中, 提高资源复用率. 实验结果表明, 本文优化后的数论变换(Number Theoretic Transform, NTT)、逆数论变换(Inverse NTT, INTT)、点对位相乘(Point-Wise Multiplication, PWM)的效率达到200 MHz, 合并执行Kyber效率达到175 MHz, 优于其他方案, 具有良好的性能.

关键词: 后量子密码; Crystal-Kyber; K²RED; 蝶形运算; 多项式乘法; 硬件效率

基金项目: 国家自然科学基金(No.61702518)

中图分类号: TP313 **文献标识码:** A **文章编号:** 0372-2112(2024)05-1679-11

电子学报URL: <http://www.ejournal.org.cn> **DOI:** 10.12263/DZXB.20220523

FPGA Efficient Parallel Optimization of Crystal-Kyber

LÜ Shun-sen^{1,2}, LI Bin^{1*}, ZHAI Jia-qi^{1,2}, LI Song-qi^{1,2}, ZHOU Qing-lei¹

(1. School of Computer and Artificial Intelligence, Zhengzhou University, Zhengzhou, Henan 450001, China;

2. State Key Laboratory of Digital Engineering and Advanced Computing, Zhengzhou, Henan 450001, China)

Abstract: Polynomial multiplication operations limit the practical applications of lattice-based post-quantum cryptography. In order to improve the performance and efficiency of post-quantum cryptography Crystal_Kyber algorithm, and reduce the running time and reduce the influence of polynomial multiplication, this paper designs a new butterfly operation unit to optimize the Kyber scheme with prime modulus $q=3329$. First of all, the algorithm is executed by 16-way parallel scheduling of the new butterfly operation unit, which shortens the calculation cycle. Secondly, using pipeline technology and improved K²RED algorithm, the design and implementation of a new butterfly operation unit for reducing resource consumption. Ultimately, the data is stored in the way of multi-RAM, and the multi-channel RAM is optimized to allow data to be stored alternately in RAM and improve the resource reuse rate. The experimental results show that the optimized NTT (number theoretic transform), INTT (Inverse NTT), PWM (point-wise multiplication) efficiency reaches 200 MHz, and the combined execution Kyber efficiency reaches 175 MHz, which is superior to other schemes and has good performance.

Key words: post quantum cryptography; Crystal-Kyber; K²RED; butterfly arithmetic; polynomial multiplication; hardware efficiency

Foundation Item(s): National Natural Science Foundation of China (No.61702518)

1 引言

随着量子计算机的高速发展, 当前传统公钥密码密体制如AES和MD5等受到了严重的挑战. 比如Shor算法能够通过量子计算机在多项式时间内破解大整数分解问题和离散对数问题. 在此背景下, 使用经典计算

机抵挡量子计算机攻击的后量子密码得到了学术界和工业界的广泛研究. 在2016年, 美国国家标准与技术研究院(National Institute of Standards and Technology, NIST)向全球发起了一场后量子密码方案竞赛, 旨在在新的量子时代到来之前形成新的公钥密码标准^[1]. 基

于格的密码算法 Crystal-Kyber^[2] 是一个用非对称加密的思想去封装密钥并进行密钥协商的算法, Kyber 支持的安全类别 1、3 和 5 分别相当于 AES-128, AES-192 和 AES-256 的强度, 它是 NIST 在经过三轮筛选后得到的后量子密码方案中四个决赛方案之一, 得到了学术界的广泛关注.

Kyber 算法由于存在大量的多项式乘法的运算, 不仅严重消耗了执行时间和资源, 而且也制约了 Kyber 算法本身的执行效率. 在现有的高效方法中, 使用数论变换 (Number Theoretic Transform, NTT) 能够加速晶格密码系统中的核心算法运算. 可以利用硬件 FPGA 来实现基于蝶形操作的 NTT 从而达到降低多项式乘法的计算复杂度^[3]的目的. 因此, 如何利用 FPGA 对 NTT 进行优化, 对 Kyber 算法的性能有着关键的作用.

受到硬件的执行效率和资源消耗的影响, 为了达到运行速度和资源的平衡, 同时保证算法本身的性能, 本文利用 FPGA 采用 16 路并行的方式对 Kyber 算法执行优化操作, 并利用改进的 K²RED 算法对 NTT, 逆数论变换 (Inverse NTT, INTT), 点对位相乘 (Point-Wise Multiplication, PWM) 中存在的蝶形运算操作进行优化, 实现了对 Kyber 算法中存在的大量多项式乘法优化的目的.

本文的贡献如下:

(1) 对 K²RED 算法进行改进, 实现一种新型蝶形运算单元, 确保高频率的同时降低了硬件资源消耗, 达到了功耗平衡的目的.

(2) 采用 16 路并行方式对蝶形单元进行调度, 提高性能的同时降低时延.

(3) 对 PWM 调度进行设计, 使新型蝶形运算单元在支持 NTT 和 INTT 运算的基础上, 支持 PWM 运算.

(4) 优化多 RAM, 在有限的空间资源下完成对 Kyber 多项式乘法的多路并行运算, 提高了空间资源利用率.

2 相关工作

本节介绍了贯穿全文的符号, 简单介绍了 Kyber 算法, 描绘了 NTT、INTT、PWM 算法的具体流程.

2.1 符号定义

整数环用 \mathbf{Z} 表示, 设置整数环模 q 的商环 $\{0, 1, \dots, q-1\}$ 使用 \mathbf{Z}_q 表示. $\mathcal{R}_q = \frac{\mathbf{Z}_q[x]}{X_n+1}$ 是定义在 \mathbf{Z}_q 为 $\mathbf{Z}/q\mathbf{Z}$ 域中整系数多项式环. 其中, X_n+1 表示的最大维度数是 $n-1$.

Kyber 算法是基于错误学习问题 (Learning With Error, LWE) 的算法^[4], 后量子安全等级见表 1. 在 NIST 第三轮公布的标准中, Kyber 的素模 q 的值由第一轮的 7 681 减少到了 3 329. 由于 Kyber 中大量的多项式乘法, 如果

采用公式法计算多项式, 会导致运算时间的增加, 所以在 Kyber 官方提供文档中采用了 NTT、INTT、PWM 变换算法来简化多项式计算.

表 1 Kyber 算法参数与安全等级

算法	NIST 等级	n	k	q
Kyber512	1	256	2	3 329
Kyber768	3	256	3	3 329
Kyber1024	5	256	4	3 329

NTT 是在快速数论变换 (Fast Fourier Transform, FFT) 数论基础上得到的一种能够对多项式乘法进行加速优化的算法, 其是在整数环 \mathbf{Z}_q 上进行的数论变换. 多项式的 NTT 运算是将一些特殊值代入多项式中, 使多项式由系数表示法转换为点值表示法. 对多项式 (1), a_i 在经过 NTT 变化后得到的新的多项式如式 (2) 所示:

$$A(x) = \sum_{i=0}^{n-1} a_i x^i, a_i \in \mathbf{Z}_q \quad (1)$$

$$\bar{a}_k = \sum_{i=0}^{n-1} a_i \omega^{ik} \bmod q \quad (2)$$

其中 ω 被称为旋转因子, 是模 q 的 n 次本原单位根, 且 $\omega \in \mathbf{Z}_q$, 满足 $\omega^n \equiv 1 \pmod{q}$, $\omega^i \neq 1 \pmod{q} \forall i < n$; 素模 q 形式一般为 $q = \rho_{\text{odd}} \times 2^{\text{pow}} + 1$, 其中 $2^{\text{pow}} > \rho_{\text{odd}}$, 满足条件 $q \equiv 1 \pmod{2n}$; INTT 的操作和 NTT 类似, 但在 INTT 中, 需要利用 ω^{-1} 来代替 ω 在 NTT 中的作用, 而且还需乘以额外的 n^{-1} . INTT 将 \bar{a}_k 转化为 a_i 如式 (3) 所示:

$$a_i = n^{-1} \sum_{k=0}^{n-1} \bar{a}_k \omega^{-ik} \bmod q \quad (3)$$

经过 NTT 之后, 多项式 A 和 B 乘积表示为 $C = \text{INTT}((\text{NTT}(A) \circ (\text{NTT}(B)))$, \circ 表示多项式的点对位相乘 PWM.

算法 1、算法 2、算法 3 给出了 NTT、INTT、PWM 在 Kyber 算法中的具体迭代流程. $\omega=17$ 是素模 $q=3 329$ 在 $n=256$ 时的第一个本原单位根, $\text{br}_{t-1}(f)$ 表示对 f 进行 $t-1$ 位的比特逆序操作.

在官方文档说明中, 利用 NTT 和 INTT 可以将原有的 256 项多项式拆分为两个独立计算的 128 项多项式. 在 PWM 算法中, 只需要进行 128 次二次多项式的乘法运算即可, 相较于传统的二项式乘法计算所需的时间更短. 但在算法 1 和算法 2 中 NTT 和 INTT 存在着大量的模加、模减和模乘. 算法 3 中, PWM 的模加和模乘运算相较于 NTT 和 INTT 更为复杂. 在硬件设计时, 大量的取模运算必然会出现繁杂的除法操作, 不仅会加剧硬件资源消耗, 而且会导致算法整体运算时间的上升. 为了提高算法效率, 需要对取模操作进行化简优化, 一般是采用其他方法替换该操作.

算法 1 NTT

输入: $A(x) = \sum_{i=0}^{n-1} a_i x^i$, $a_i = \{a_0, a_1, \dots, a_{n-1}\}$, ω, q

输出: $\bar{A}(x) = \sum_{i=0}^{n-1} \bar{a}_i x^i$, $\bar{a}_i = \{\bar{a}_0, \bar{a}_1, \dots, \bar{a}_{n-1}\}$

1. FOR $i = (\log_2 n) - 1$ TO 1 DO
2. $m = 2^i$ $r = 0$
3. FOR $k = 0$ TO $n - 1$ BY m DO
4. $W = \left(\omega^{\text{br}_{(\log_2 n - i)(r)}} \right) \bmod q$
5. FOR $j = 0$ TO $m/2 - 1$ DO
6. $u = a_{j+k}$ $v = a_{j+k+m/2}$
7. $t = (W \times v) \bmod q$
8. $e = (u + t) \bmod q$ $o = (v - t) \bmod q$
9. $a_{j+k} = e$ $a_{j+k+m/2} = o$
10. END FOR
11. $r = r + 1$
12. END FOR
13. END FOR
14. RETURN $\bar{A}(x)$

算法 2 INTT

输入: $\bar{A}(x) = \sum_{i=0}^{n-1} \bar{a}_i x^i$, $\bar{a}_i = \{\bar{a}_0, \bar{a}_1, \dots, \bar{a}_{n-1}\}$

输出: $A(x) = \sum_{i=0}^{n-1} a_i x^i$, $a_i = \{a_0, a_1, \dots, a_{n-1}\}$

- 1 $k = 0$
2. FOR $i = (\log_2 n - 1)$ TO 1 DO
3. $m = 2^{(\log_2 n - i)}$
4. FOR $s = 0$ BY m TO 2^i DO
5. $u = \bar{a}_j$, $v = \bar{a}_{j+k}$, $W = \omega^{\text{br}_{(\log_2 n - i)(k) + 1}} \bmod q$
6. $e = (u + v) \bmod q$, $o = (u - v)W \bmod q$
7. $\bar{a}_j = (e/2) \bmod q$, $\bar{a}_{j+k} = (o/2) \bmod q$
8. END FOR
9. $k = k + 1$
10. END FOR
11. RETURN $A(x)$

算法 3 PWM

输入: $\bar{A}(x) = \sum_{i=0}^{n-1} \bar{a}_i x^i$, $\bar{a}_i = \{\bar{a}_0, \bar{a}_1, \dots, \bar{a}_{n-1}\}$

$\bar{B}(x) = \sum_{i=0}^{n-1} \bar{b}_i x^i$, $\bar{b}_i = \{\bar{b}_0, \bar{b}_1, \dots, \bar{b}_{n-1}\}$

输出: $\bar{C}(x) = \sum_{i=0}^{n-1} \bar{c}_i x^i$, $a_i = \{\bar{c}_0, \bar{c}_1, \dots, \bar{c}_{n-1}\}$

1. FOR i FROM 0 TO $2^{(\log_2 n) - 1}$ DO
2. $W = \left(\omega^{1 + \text{br}_{\log_2 n - (i)}} \right) \bmod q$
3. $a_0 = \bar{a}[2i]$, $a_1 = \bar{a}[2i + 1]$
4. $b_0 = \bar{b}[2i]$, $b_1 = \bar{b}[2i + 1]$
5. $\bar{c}[2i] = (a_0 b_1 + a_1 b_0) \bmod q$
6. $\bar{c}[2i + 1] = (a_1 b_1 W + a_0 b_0) \bmod q$
7. END FOR
8. RETURN $\bar{C}(x)$

等^[7]使用了集成了有限域指令集扩展的 RISC-V 处理器来加速 Kyber 轻量级架构中的多项式乘法,相较于软件实现提速 85%,但是 Alkim 等使用处理器的方式优化算法,资源消耗很大,得到的频率较差,在实际应用中的场景受限. Huang 等^[8]、Dang 等^[9]、Xing 等^[10]提出了 Kyber 的纯硬件结构体系,但是 Huang 等的工作依赖组件之间的 Block RAM 原语来执行算术任务和存储中间结果, Xing 等采用基于时域和频率的结合型蝶形单元来实现 NTT 操作,虽然在面积和时间上有了很好的提升,但是频率还可以进一步优化. Chen 等^[11]提出一种用于计算 NTT 和点对位乘法的 Kyber 处理器,达到了加快 NTT 计算频率目的,但并没有和算法性能达到平衡,消耗了过多资源. Roy 等^[12]利用传统的二次复杂度多项式乘法算法设计了一种灵活的硬件结构,但该硬件结构加重了资源负担,没有对 Kyber 算法起到很好的优化作用.

国内对后量子密码的研究和国外相比较晚,但发展的很快. 周朕等^[13]针对国内研发的 Aegis-enc^[14]设计了一种基于 NTT 算法实现环上多项式运算的紧凑硬件架构,该硬件架构采用 Montgomery 算法对关键模乘部分进行优化,算法整体性能得到了提高. 崔益军等^[15]针对 NTT 优化设计了迭代型和流水线硬件结构,他们采用了文献^[10]提出的修改的 Barrett 约减算法对 NTT 的关键部分进行优化,达到了提高硬件的执行效率的目的. Feng 等利用 Stockham FFT 算法设计了一种多路径且高并行度的 NTT 算法^[16],并且通过利用大量蝶形单元并行调度提高 NTT 模块的运行速度.

综上分析可知,利用 FPGA 对 Kyber 算法进行优化时,为了提高算法的性能. 需要对 Kyber 算法中的关键部分操作进行深度优化,化繁为简;而如果要提高算法的运行速度,增加并行度的同时提高各模块之间的数

2.2 国内外研究现状

Botros 等^[5]研究了 Kyber 算法的软件实现,提出了在 Cortex-M4 上对 Kyber 算法进行了优化,但是采用软件方式实现 Kyber 核心运算,不仅消耗时间,而且得到的频率较低,无法满足高性能的需求. 也有诸多学者提出使用硬件的方式来对 Kyber 算法进行优化. Basu 等^[6]基于 HLS 的硬件设计方法将 Kyber 算法映射到 FPGA 上进行实现,和软件实现相比,该方法性能有了很大的提升,但是 Basu 等人只是对算法进行比较,并没有对算法进行深度优化,使用的硬件资源过大,频率较低. Alkim

据处理能力是一个很好的选择. 但大部分方案中在对 NTT 进行优化的时候, 由于硬件资源问题, 舍弃了高并行, 导致 NTT 计算周期长. 而有些方案在采用高并行的方式加速算法时, 也存在消耗过多的硬件资源的问题. 且部分方案为了进一步提高算法性能仅执行了 NTT 和 INTT 运算, 舍弃了 PWM 运算. 因此如何做到满足算法性能的同时与硬件资源达到平衡, 是当前学术界在格密码方案硬件实现的一大研究热点.

本文从 NTT、INTT、PWM 各计算阶段着手, 设计实现一种新型蝶形单元, 使得该蝶形单元支持 NTT、INTT、PWM 运算, 实现了对 Kyber 算法关键部分的化简, 并采用高度并行方式进行调度, 提高计算速度的同时实现了功耗平衡.

3 Kyber 优化设计

3.1 整体架构

为在提高 Kyber 算法运算速度的同时, 增强硬件整体的灵活性和可扩展性, 本文首先对 NTT、PWM、INTT 和 Kyber 多项式乘法进行 16 路并行执行, 设置了 16 个蝶形单元匹配执行运算操作. 其次, 采用多 RAM 的方式, 对执行运算的多项式参数和每轮蝶形操作迭代后的结果进行存取; 对于旋转因子 ω , 采用预训练的方式提前计算好之后保存在 ROM 中. 最后, 对旋转因子 ω 的读取控制、RAM 存取地址的索引控制和蝶形单元的输入输出都由控制逻辑来完成. 图 1 显示了 FPGA 整体结构.

如图 1 所示, 16 路蝶形运算操作并行执行, 蝶形单元内部采用流水线技术进行实现; 为了同步 RAM 的存取地址索引, 在控制单元内将 Kyber 执行蝶形单元进行的轮数与循环次数和 RAM 的存储空间地址进行绑定, 减少了程序的复杂性. 在 PWM 阶段, 因为 K^2RED 算法的特殊性, 蝶形单元输出的最终结果由 T 值代替了 O 值. 为了匹配 16 路蝶形运算, 需要在每轮蝶形运算时从 ROM 中取出该轮蝶形运算对应的 16 个旋转因子. 硬件整体采用松耦合架构展开, 可以根据输入的控制信号的不同, 自主独立完成 NTT、INTT、PWM 运算.

3.2 流水线型蝶形单元硬件设计

16 路蝶形单元块同步进行运算, 虽然显著提升了运算速度, 但也导致了硬件资源消耗过大. 为进一步节省资源, 首先, 在蝶形单元内部只设计了一个模加器和模减器, 并采用组合逻辑的方法实现模加减. 其次, 利用改进的 K^2RED 算法对 Kyber 中耗时的求模运算进行化简. 由于需要在使用改进的 K^2RED 算法前对传递进来的参数进行相乘, 考虑到 DSP 块在最新的 FPGA 设备上大量可用, 为了充分利用 FPGA 资源, 选择使用 DSP 作为基础乘法器, 从而实现快速高效的系数级乘法计算, 进一步降低计算时间. 最后, 因为 NTT、INTT、PWM

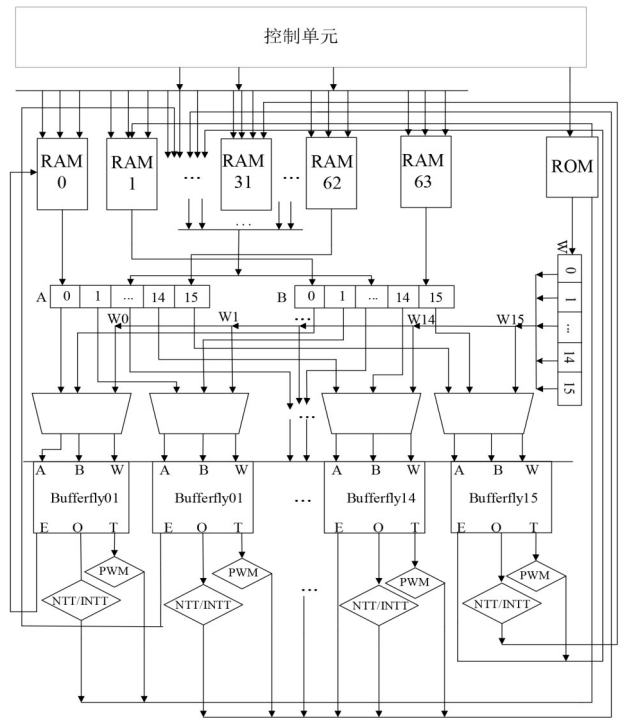


图 1 FPGA 整体结构

算法执行流程的不同, 在不同的阶段内, 保存在 RAM 块中的结果也不同, 为了减少顶层逻辑的复杂度, 以及进一步降低路径延时: 第一, 在蝶形运算单元内部设计了多个缓存器来缓存中间结果; 第二, 在顶层设置仲裁信号, 根据不同的信号执行不同的流程. 该蝶形单元硬件结构和其各阶段流水线执行流程图如图 2 和图 3 所示.

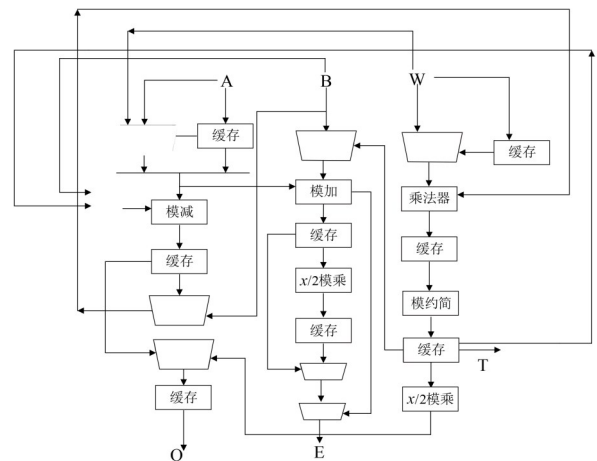


图 2 蝶形运算单元硬件结构

在图 2 中, 蝶形单元内部只设置了一个模加器、模减器和模约简器. 它还有用于同步输出系数的缓存寄存器 and 用于可重构的多路复用器. 该蝶形单元以 A、B、W 为输入进行运算操作, 在对应的阶段内生成结果值 E、O(T).

如图 3 所示,在蝶形单元内部,NTT、INTT、PWM 在不同的阶段内执行不同的流程.首先,为了使算法执行流程和硬件时钟相匹配.由于改进的 K^2RED 算法的特殊性,需要对 PWM 调度进行改进.其次,为了进一步节省资源,缩短硬件的计算周期,对模加减、 $x/2 \bmod q$ 进行进一步的优化.最后,为了匹配不同的流程,在时钟、RAM 索引调度、控制单元上进行合理的设计,从而得到正确的结果.

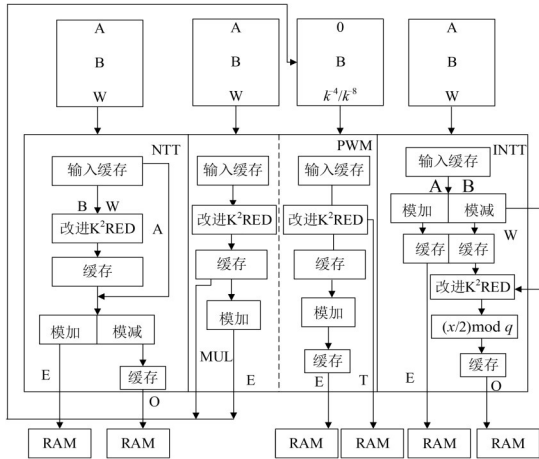


图 3 流水线执行流程

本文采用下文所述方法对 K^2RED 算法进行改进,使该算法能够支持 NTT、INTT、PWM 运算;然后对 PWM 调度进行设计,使其能够匹配改进的 K^2RED 算法;在此基础上对模加减、 $x/2 \bmod q$ 进行优化;最后,对多 RAM 存取进行设计,进一步增加空间资源复用.

3.2.1 改进的 K^2RED 算法

针对蝶形单元运算中最耗时的求模运算操作,本文采用文献[3]推荐的 K^2RED 算法对该操作进行优化.之所以使用该算法,是因为 K^2RED 算法对任意整数 C 的求模运算转化为了对 C 的部分 bit 拆分移位操作得到整数 D ,使 $D = k^2 C \bmod q$. 在该流程中,除法运算被舍弃,而且硬件中的 bit 移位操作实现简单且不消耗资源.在 Kyber 算法中,得益于 $q=3\ 329$,其形式可以表示为 $q = k \times 2^m + 1$,其中 $k=13, m=8$,基于该性质, K^2RED 算法可以在蝶形单元中被使用对求模运算进行化简.

文献[3]设计的 K^2RED 模块只支持 NTT 和 INTT 运算,且涉及了 k 的乘法操作.本文对该算法进行改进和优化,舍弃 K^2RED 算法中耗时的乘法操作,而只选择使用 bit 判断、移位、加法和减法,进一步提高了算法在硬件中的执行速度,降低了资源消耗,并使其在支持 NTT 和 INTT 运算的同时也支持 PWM 运算.改进的 K^2RED 算法的具体流程如算法 4 所示.

改进的 K^2RED 算法利用输入数本身的 bit 移位操作代替高复杂度的除法来实现求模运算,在节省硬件

算法 4 改进的 K^2RED 算法

输入: $C = A \times B$

输出: $R = k^2(A \times B) \bmod q$

1. $c_1 = C[7:0]; c_2 = C[23:8]$
2. $c_3 = (c_1 \ll 3) + (c_1 \ll 2) + c_1 - c_2$
3. $c_4 = c_3[7:0]; c_5 = c_3[15:8]$
4. $c_6 = (c_3 \ll 3) + (c_4 \ll 2) + c_3 - c_4$
5. $c_7 = (c_5[7:5] = 0) ? \#c_6 : c_6 + 256$
6. $c_8 = c_7 - 3329$
7. $R = (c_7[12] = 0) ? \#((c_8[12] = 0) ? (c_8 : c_7)) : (c_7 + 3329)$
8. RETURN R

资源的同时也缩短了计算周期.

3.2.2 PWM 调度设计

在设置 PWM 模块硬件结构时,由 PWM 算法可知,需要计算 $a_0 b_1, a_1 b_0, a_1 b_1 W, a_0 b_0$,共 5 次模乘和 2 次模加.而且由算法 3 可知, PWM 算法流程的第 3~6 步和 NTT、INTT 类似,所以都可以采用 K^2RED 算法对求模操作进行化简.但是由于 K^2RED 算法返回值 $D = k^2 C \bmod q$,需要消除 k^2 所带来的影响后才能进行相乘操作,从而得到正确的 $\bar{c}[2i]$ 和 $\bar{c}[2i+1]$.

在对 PWM 调度进行设计的时候,我们先在前 5 个时钟完成 PWM 算法的第 3~6 步,得到被 k^2 影响的结果值,在后续 5 个时钟内消除 k^2 对结果值所带来的影响.经过研究发现,需要在特定的时钟内将蝶形单元输出值与模逆值 $k^{-4} = 0x549, k^{-8} = 0x8f5$ 进行结合后再次进行蝶形操作.一共需要 10 个时钟,图 4 给出了 PWM 阶段具体调度流程.

如图 4 所示,在前 5 个时钟完成 PWM 算法的 3~6 的流程,然后在后 5 个时钟消除 k^2 的影响,其中在第 9 个时钟的时候,因为需要使用的 ADD 值是 K^2RED 运算后经过模加后得到的,所以需要等待 1 个时钟.

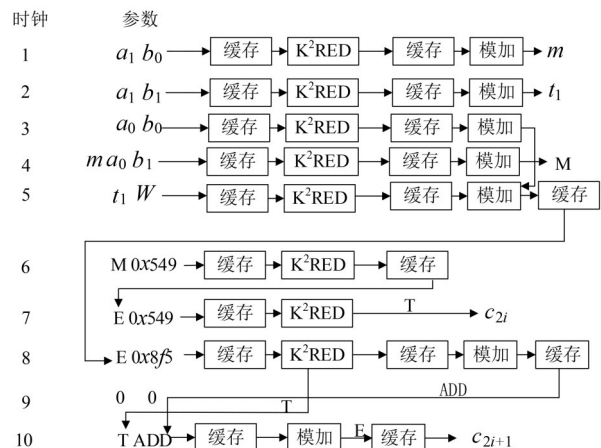


图 4 PWM 调度流程图

3.2.3 模加减与 $x/2$ 模乘

由图2可知, NTT、INTT、PWM 都在执行阶段内使用到了模加减, 蝶形单元结果值 $E、O(T)$ 和模加减息息相关. 模加减对蝶形运算起到重要作用. 一个好的模加减, 不仅能加速算法的执行时间, 而且能进一步降低蝶形运算单元的资源消耗. 本文对模加减采用组合逻辑实现来避免除法, 具体流程是判断模加减的结果值和素数 q 的关系后进行处理. 其硬件结构如图5、图6所示.

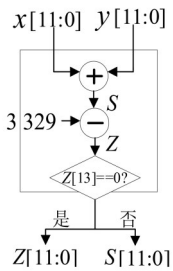


图5 模加硬件结构

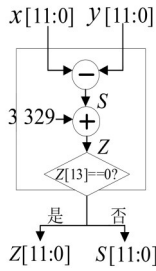


图6 模减硬件结构

对于 $x/2 \bmod q$ 操作, 因为涉及取模和乘 $1/2$ 操作, 本文采用了 Zhang 等^[17]提出的方法对 $x/2 \bmod q$ 进行优化. 具体是由于 q 为常数 3 329 且是奇素数的性质, 对于 $x/2 \bmod q$ 操作, 可以通过移位、加法和逻辑与操作简化乘法和取模预算, 进一步节省资源开销, 其计算式如式(4)所示. 根据公式所设计的硬件结构如图7所示, G 表示输入数 x 最低比特位 $x[0]$ 的 10 倍扩展.

$$\frac{x}{2} \bmod q = (x \gg 1) + x[0] \& \left(\frac{q+1}{2} \right) \quad (4)$$

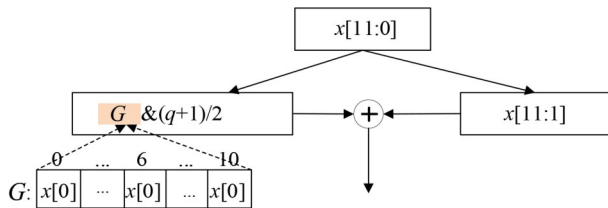


图7 $x/2 \bmod$ 优化结构图

3.3 多 RAM 存取设计

在 Kyber 算法流程中, 需要将多项式 A 和 B 的 NTT 运算结果值执行 PWM 运算, 然后对 PWM 结果执行 INTT 运算. 考虑到对于 N 个多项式点, 需要执行 $\log_2 N$ 的 NTT 和 INTT, 为了便于后续操作, 需要将每轮蝶形单元运算得到的结果保存在 RAM 中. 在该过程中, 因为涉及多个多项式的参数存取, 而且需要防止执行蝶形操作的参数混乱, 所以对保存结果的 RAM 设置唯一的存取地址.

图8所示的8点多项式蝶形运算, 既要保证每个阶段内进行蝶形运算取得的多项式参数正确, 也要确保蝶形运算中间结果必须存储在指定的索引位置上, 从而防止索引地址混淆. 而在 NIST 第三轮的 Kyber 算法中, 进行蝶形运算的每个多项式参数有 256 个, 因为 NTT 和 INTT 的原因, 可以将一个 256 项的多项式拆分为两个独立计算的 128 次多项式进行蝶形运算, 一共需要进行 $\log_2 256 = 8$ 轮蝶形运算才能最终结果. 对多 RAM 的索引控制更为复杂, 所需要的 RAM 块更多. 为了进一步优化 FPGA 的空间资源, 减少资源消耗, 以下给出本文对多 RAM 的优化方法.

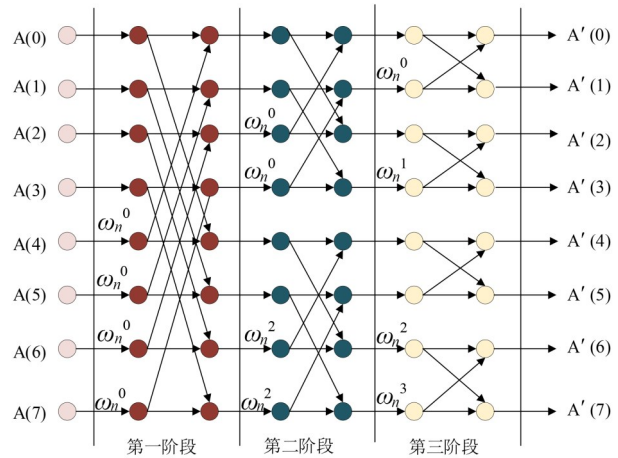


图8 8点多项式蝶形运算

3.3.1 多 RAM 优化

RAM 是 FPGA 提供的紧凑的真正双端口内存库, 方便存储大量数据. 在本文中为了在一个时钟内执行 16 次蝶形运算, 一共设计了 16×4 个 RAM 块用来存储, 单个 RAM 块需要有 $(256 \times 2) / 64 = 8$ 个存储空间来存储多项式参数. 前 32 个 RAM 块存储第一个多项式的参数, 后 32 个 RAM 块存储第二个多项式参数. 在执行 NTT 和 INTT 的时候, 在同一时钟内分别从各自 32 个 RAM 块中选择取出执行蝶形操作的第 $2 \times i$ 和 $2 \times i + 1$ ($0 \leq i \leq 16$) 的多项式参数进行 16 路并行执行 NTT 和 INTT. 如此, 16 个蝶形运算单元可独立并行存取参数并执行. 图9给出了 NTT 阶段首轮初始化参数存储 RAM 示意图. 如图9所示, 将 256 项多项式平分为上下 128 项, 利用偶数组 RAM 存储上 128 项, 奇数组 RAM 存储下 128 项. NTT 首轮循环开始后, 从第一部分偶数组中取出 a_0, a_1, \dots, a_{15} 与第一部分奇数组 $a_{128}, a_{129}, \dots, a_{143}$ 分别组合进行蝶形运算. 随后, 对第二部分按照相同方式组合进行蝶形运算, 直到第一轮 NTT 结束. 同组相同部分之间多项式索引地址连续. 单 RAM 内相邻部分多项式索引地址相差 16.

NTT、INTT、PWM 合并执行 Kyber 多项式乘法运算

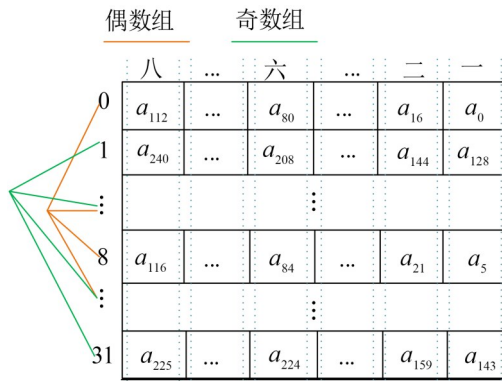


图9 NTT初始参数存储

时,单RAM存储示意图如图10所示.对于多RAM,将64个RAM块平分分成两组,每组按照索引地址平分分成集合1和集合2.将初始化要执行Kyber多项式乘法运算的多项式 $A(x)$ 和 $B(x)$ 分别保存在两个RAM分组的集合1中,保存完毕后从集合1取出多项式参数执行NTT,将NTT运算结果保存在集合2中;在PWM阶段取出集合2中元素进行蝶形运算,并将结果保存在集合1中,等待INTT阶段的执行.多RAM存储流程如图11所示.

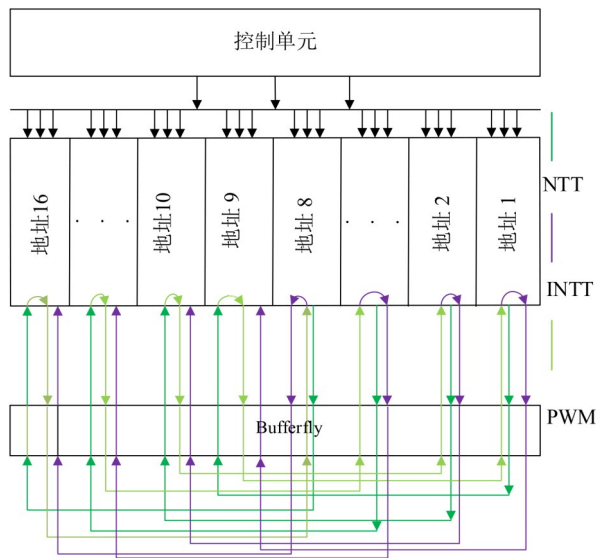


图10 单RAM存储图

最后,由于NTT、INTT、PWM在执行相应轮数时对旋转因子 ω 值的参数固定,因此本文在对旋转因子 ω 进行预处理后,按照一定的顺序将其存放在BROM中.具体来说,地址0~38存放执行NTT的旋转因子 ω .地址39~77存放执行INTT的旋转因子 ω' .地址78~85存放执行PWM的选择因子 ω^* .在每个地址存放16个旋转因子,对应并行执行的16个蝶形运算单元.ROM存储结构如图12所示.

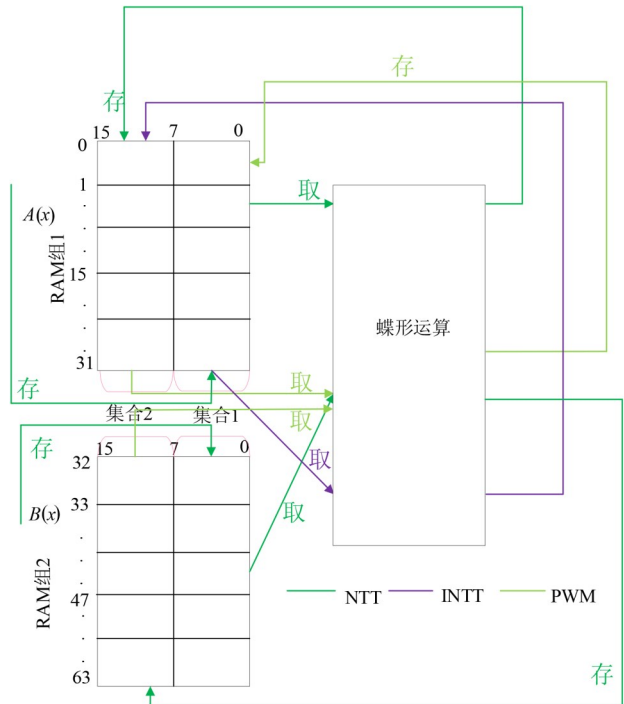


图11 多RAM存储流程图

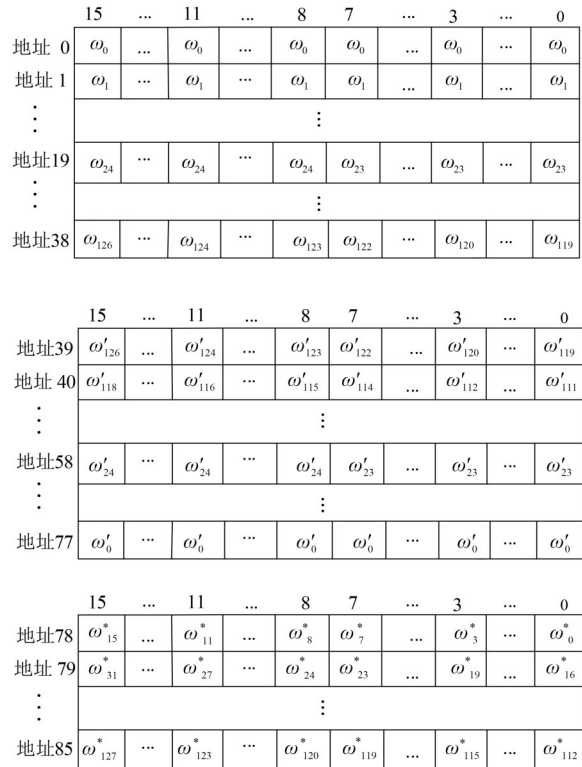


图12 ROM存储图

3.4 RAM块索引地址控制

从第3.3节的分析可知,NTT、INTT、PWM三个阶段中每轮进行蝶形运算的多项式参数不同,但存储在RAM块中的每个多项式参数地址唯一且固定,且需要

将每轮蝶形运算执行后得到的结果 $E、O(T)$ 保存在 RAM 中指定索引位置上,以便下一个阶段读取,保证蝶形运算得到的结果值的正确性。

在 RAM 地址块设计时采用 1 个地址 red 对 RAM 块进行读取地址控制,两个地址 wre 和 wro 对 $E、O(T)$ 进行存储地址控制。在设置 rae、wre、wro 的参数时,考虑到 NTT、INTT、PWM 经过 7 轮蝶形运算完成存取,而且 RAM 块中存储空间最高最有 16 个,令 rae、wre、wro 的后三比特位由 NTT、PWM、INTT 的经历的轮数和循环次数决定。另外,由于 RAM 中高 8 位和低 8 位的存取是循环进行的。因此对 rae、wre、wro 的最高比特位按照每 7 轮进行一次交替取反。算法 5 和算法 6 介绍了 NTT 情况下的 rae、wre、wro 的具体计算流程。

算法 5 RAM 读取地址算法

输入: 蝶形运算阶段 stage[2:0], 循环次数 loop[2:0]

输出: RAM 读地址 rae[2:0]

1. 初始化 rae[2:0]=0
2. IF (stage > 1)
 - rae[2:0]=loop
3. ELSE
 - rae[2:0]=((loop >> 1) & ((1 << (2 - stage)) - 1)) + (loop[0] << (2 - stage)) + ((loop >> (3 - stage)) << (3 - stage))
4. RETURN rae[2:0]

如算法 5 所示,在步骤 2 中,当蝶形运算阶段大于 1 时,RAM 的读地址直接由循环次数决定。否则,如步骤 3 所示,读地址需要根据当前 stage 和 loop 进行额外计

算法 6 RAM 写入地址算法

输入: 蝶形运算阶段 stage[2:0], 循环次数 loop[2:0]

输出: RAM 写入地址 wre[2:0], wro[2:0]

1. 初始化 wre[2:0]=0; wro[2:0]=0
2. IF (stage > 2)
 - wre[2:0]=loop; wro[2:0]=loop
3. ELSE
 - wre[2:0]=(loop >> 1) + ((loop >> (3 - stage)) << (2 - stage));
 - wro[2:0]=(loop >> 1) + ((loop >> (3 - stage)) << (2 - stage)) + (1 << (2 - stage));
4. RETURN wre[2:0]; wro[2:0]

算,从而决定当前 RAM 的读地址。

如算法 6 中,因为 RAM 是先读后写,而且需要等待 3 个时钟。所以在步骤 2 中,只有 stage 大于 2 的时候,wre 和 wro 的值由 loop 决定。而在其他情况也是根据当前 stage 和 loop 决定。

考虑到算法 5 和算法 6 在一定阶段的特殊性,在硬件设计时,针对 RAM 读写索引启动了延时缓存策略。具体在 NTT 蝶形运算的前 3 轮,由于此时 NTT 的读写不是按序进行的,在每轮 NTT 蝶形结束运行时等待 2 个时钟,在 4~7 轮,此时 RAM 读写顺序一致,无需进行中间等待。在 INTT 模块,在前三轮蝶形运算时每轮等待 1 个时钟。对于 PWM 模块,由于 PWM 阶段每轮运算都需要经过多个状态,所以有足够的时钟使得读写地址匹配。图 13 显示了 NTT 阶段的 RAM 读写索引控制。

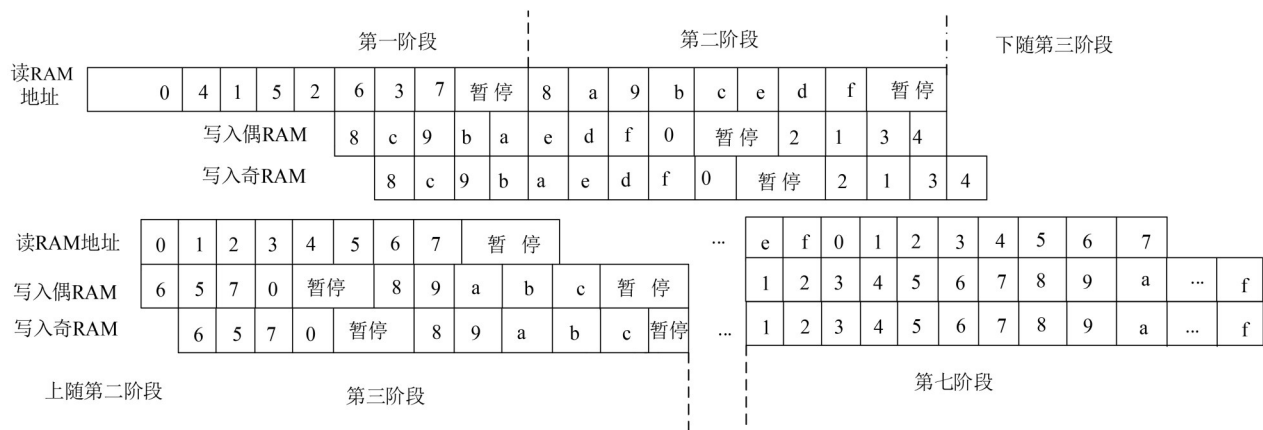


图 13 RAM 读写索引控制

为解决多蝶形单元并行计算地址冲突问题,每 2 个 RAM 对应 1 个蝶形运算,将 32 个 RAM 块分配给其对应的 16 个蝶形单元,各蝶形单元独立执行。其对应的 RAM 块内按照算法 5、算法 6 对蝶形运算系数和运算结果进行读取。

4 实验结果分析

4.1 实验结果与分析

为了验证本文所提出的优化方案,我们在 ZYNQ-XC7Z035FFG676-2 上实现了该程序。表 2 展示了 Kyber 算法中各模块的资源使用情况。表 3 给出了本方案中

NTT、INTT、PWM 和 Kyber 总体的吞吐量和资源使用情况。

表 2 Kyber 算法模块资源使用表

模块	实现方式	LUT	FF	BRAM	DSP
地址控制单元	流水线	3 351	134	0	0
蝶形运算单元	流水线	255	132	0	1
多RAM存储	并行	256	0	32	0
ROM存储块	串行	0	0	2.5	0
Kyber实现	串并混合	8 053	3 321	34.5	16

表 3 Kyber 性能表

模块	LUT	FF	BRAM	DSP	频率/MHz	功耗/W
NTT	4 981	3 063	18.5	16	200	0.690
INTT	4 418	2 797	18.5	16	200	0.661
PWM	4 470	2 790	18.5	16	200	0.670
Kyber	8 053	3 321	34.5	16	175	0.774

从表 2 分析可知,地址控制单元由于需要计算蝶形运算轮数、本轮迭代次数、RAM 和 ROM 的读写地址、控制信号以及读写等待信号等,所以消耗了更多的 LUT. 蝶形运算单元占用了较少的资源实现了模加、模减、模乘等运算,且仅需要 1 个 DSP 便可完成乘法运算,表明使用改进 K²RED 优化蝶形运算单元的正

确性. 64 个 RAM 块仅使用了很少的 LUT 和 BRAM 资源便完成了 Kyber 多项式的参数存储. BROM 块使用了少量的 BRAM 资源,便完成了对 Kyber 多项式乘法所需的旋转因子 ω 的存储. 所使用的 LUT 占用芯片 LUT 资源的 4.68%, FF 占用比例为 0.97%, BRAM 资源占用比例为 6.90%, DSP 资源占用比例为 1.78% 便完成了 Kyber 多项式乘法在 FPGA 的实现,达到了资源上的优化目的,证明了本文优化思想的正确性.

4.2 方案对比

为了更加方便的比较性能优劣,本文引入了一个新的计量值/参数/指标——单元性能 (Unit Performance, UP),客观反应了资源消耗和算法性能之间的关系. UP 值越高,表明在有限的硬件资源条件下,优化后的算法性能和资源利用率越好. UP 计算方式如下所示.

$$UP = \frac{\text{频率} \times 10^3}{FF + LUT} \quad (5)$$

在表 4 中,将本文算法与其他 Kyber 算法经过 FPGA 优化后的资源消耗、频率和 UP 等进行了综合对比. 在表 5 中,将本文算法与其他 Kyber 算法在 NTT 阶段经过 FPGA 优化后的资源消耗、频率和 UP 等进行了综合对比.

表 4 FPGA 优化后 Kyber 总体综合对比

方案	设备	LUT	FF	BRAM	DSP	频率/MHz	时间/ μ s	UP
文献[9]	Artix-7	12 200	12 400	15	8	210	—	8.54
文献[18]	Alveo U250	166 100	46 000	16	138	—	1.97	—
文献[19]	Artix-7	9 508	2 684	35	16	172	0.40	14.1
本文	Zynq-7035	8 053	3 321	34.5	16	175	1.50	15.38

表 4 比较了不同方案优化后的 FPGA 实现的 Kyber 多项式乘法的优化性能. 本文通过改进的 K²RED 算法优化蝶形操作,并且利用多路并行的方式优化后的 Kyber 算法在较高的频率下使用的资源相比其他文献使用的更少,得到的 UP 更好,资源利用率更好. 在表 4 中,文献[9]虽然得到的频率比本文方案多,但是其消耗的资源更多,而且从 UP 方面的比较来看,本文对资源使用和算法性能的平衡更好. 文献[18]采用 RISC-V 矢量扩展来加速 Kyber 多项式乘法,且更关注算法的执行时间而忽略了硬件资源消耗,增加了电路的复杂度. 在文献[19]中,采用 16 个蝶形运算单元的方式来实现 NTT,但是其采用的模约简方式相较于本文更复杂,导致消耗的资源更多.

在表 5 中,为了进一步比较优化后的算法性能,我们对 Kyber 算法在 FPGA 上实现 NTT 进行了性能比较. 文献[3]采用原版的 K²RED 算法对 NTT 进行优化,但文献[3]采用了时间换取资源的策略,使用了一个 2×2 蝶

型核搭配一个 K²RED 模块按照完全流水线的方式执行 NTT,相较于本文的 16 个蝶形单元使用的资源必定很少,频率更高. 但其优化后得到的 NTT 吞吐量却和本文近似相同. 在计算周期上,本文可在 0.40 μ s 内完成对 NTT 的计算,文献[3]却需要 1.46 μ s,比本文所需的时间更长. 综合对比下,本文的性能更好. 文献[20]采用 K-RED 算法对求模运算进行化简并搭配 2×2 蝶型核对 NTT 进行优化,虽然降低了硬件资源消耗,但该方案的频率、执行时间都不如本方案. 文献[21]采用 Barrett 算法对蝶形单元的求模运算进行优化,并放置 32 个蝶形单元执行 NTT. 从表 5 的分析可知,本文利用改进的 K²RED 算法优化后的蝶形单元的性能相比文献[21]的方案更好,UP 值更高. 文献[22]在 FPGA 内放置 4 个蝶形单元和 8 个乘法器并采用 SAMS2 算法^[22]对求模运算进行化简实现 NTT,在综合性能和 UP 等方面都不如本文提出的方案. 文献[23]和文献[5]同样采用 Barrett 模约简来对蝶形运算单元进行优化,但其采用乒乓操作

表5 FPGA优化后NTT总体综合对比

方案	设备	LUT	FF	BRAM	DSP	频率/MHz	时间/ μ s	UP
文献[3]	Artix-7	801	717	2	4	222	1.46	146.24
文献[20]	Zynq-7000	2 832	1 381	10	8	150	17.44	35.6
文献[21]	Zynq-7035	8 428	3 979	11	32	175	0.25	14.1
文献[22]	Artix-7	4 823	2 901	10	8	153	—	19.8
文献[23]	Artix-7	18 000	5 000	15	6	115	14.8	5.00
本文	Zynq-7035	4 981	3 063	18.5	16	200	0.40	24.86

对系数存取进行优化,虽然提高了NTT的计算速度,但是整体频率偏低,消耗资源过高.而本方案中实现的NTT在得到高频率的同时也兼顾了资源消耗.

对比Kyber和NTT的整体性能之后可知,本文提出的Kyber的16路并行优化方案在保证较高的工作频率的同时,也充分利用了芯片的资源,效率更好,具有更优的使用价值.

5 结束语

随着计算机的快速发展,基于格的密码方案在即将到来的量子计算机时代具有优异的研究价值和实用场景.本文利用改进的 K^2 RED算法对蝶形运算操作进行优化,并采用流水线技术在硬件上对优化后的蝶形单元进行实现,在此基础上通过16路并行方式完成了对NTT、INTT、PWM以及Kyber多项式的运算操作.为了进一步优化空间资源的利用率,本文采用多RAM的方式对多项式进行了存取优化.实验结果证明,本文的实验方案具有较高的实际应用和研究价值.

在未来工作中,将本文对Kyber算法的优化应用到其他格密码方案,是我们研究的重点之一.并且我们将对PWM阶段继续进行优化,进一步缩短蝶形运算时钟周期.

参考文献

- [1] NIST. Post-Quantum cryptography call for proposals[EB/OL](2017-01-03)[2022-04-15]. <https://csrc.nist.gov/Projects/post-quantum-cryptography/post-quantum-cryptograph-standards/Call-for-Proposals>.
- [2] AVANZI R, BOS J, DUCAS L, et al. Crystals-kyber[R]. Maryland: NIST, 2017.
- [3] BISHEH-NIASAR M, AZARDERAKHSH R, MOZAFARI-KERMANI M. High-speed NTT-based polynomial multiplication accelerator for post-quantum cryptography [C]//2021 IEEE 28th Symposium on Computer Arithmetic (ARITH). Piscataway:IEEE, 2021: 94-101.
- [4] REGEV O. On lattices, learning with errors, random linear codes, and cryptography[J]. Journal of the ACM, 2009, 56(6): 1-40.
- [5] BOTROS L, KANNWISCHER M J, SCHWABE P. Memory-efficient high-speed implementation of Kyber on Cortex-M4[C]//International Conference on Cryptology in Africa. Springer: Cham, 2019: 209-228.
- [6] BASU K, SONI D, NABEEL M, et al. Nist post-quantum cryptography-a hardware evaluation study[J]. IACR Cryptology Eprint Archive, 2019, 47: 1-16.
- [7] ALKIM E, EVKAN H, LAHR N, et al. ISA extensions for finite field arithmetic: Accelerating Kyber and NewHope on RISC-V[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2020, 2020: 219-242.
- [8] HUANG Y M, HUANG M Q, LEI Z K, et al. A pure hardware implementation of Crystals-Kyber PQC algorithm through resource reuse[J]. IEICE Electronics Express, 2020, 17(17): 1-6.
- [9] DANG V, FARAHMAND F, ANDRZEJ-CZAK M, et al. Implementation and benchmarking of round 2 candidates in the NIST post-quantum cryptography standardization process using hardware and software/hardware co-design approaches[EB/OL]. (2020)[2022]. <https://eprint.iacr.org/2020/795>.
- [10] XING Y F, LI S G. A compact hardware implementation of CCA-secure key exchange mechanism Crystals-Kyber on FPGA[J]. IACR Transactions on Cryptographic Hardware and Embedded Systems, 2021, 2021(2): 328-356.
- [11] CHEN Z H, MA Y, CHEN T Y, et al. Towards efficient Kyber on FPGAs: A processor for vector of polynomials [C]//Proceedings of 2020 25th Asia and South Pacific Design Automation Conference(ASP-DAC). Piscataway: IEEE Press, 2020: 247-252.
- [12] ROY S 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.
- [13] 周朕,何德彪,罗敏,等.紧凑的Aigis-sig数字签名方案软硬件协同实现方法[J].网络与信息安全学报,2021,7(2): 64-76.

ZHOU Z, HE D B, LUO M, et al. Compact software/hard-

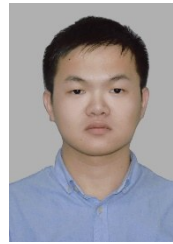
ware co-design and implementation method of Aegis-sig digital signature scheme[J]. Chinese Journal of Network and Information Security, 2021, 7(2): 64-76. (in Chinese)

- [14] ZHANG J, YU Y, FAN S, et al. Tweaking the asymmetry of asymmetric-key cryptography on lattices: KEMs and signatures of smaller sizes[C]//IACR International Conference on Public-Key Cryptography. Cham: Springer, 2020: 37-65.
- [15] 崔益军, 姚衍, 倪子颖, 等. 基于 MLWE 的格密码高效硬件实现[J]. 信息安全学报, 2021, 6(6): 40-50.
CUI Y J, YAO K, NI Z Y, et al. Efficient hardware implementation of MLWE lattice based cryptography[J]. Journal of Cyber Security, 2021, 6(6): 40-50. (in Chinese)
- [16] FENG X, LI S, XU S. RLWE-oriented high-speed polynomial multiplier utilizing multi-lane stockham NTT algorithm[J]. IEEE Transactions on Circuits and Systems II: Express Briefs, 2019, 67(3): 556-559.
- [17] ZHANG N, YANG B, 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, 2: 49-72.
- [18] LI H M, MENTENS N, PICEK S. A scalable SIMD RISC-V based processor with customized vector extensions for CRYSTALS-kyber[C]//Proceedings of the 59th ACM/IEEE Design Automation Conference. New York: ACM, 2022: 733-738.
- [19] YAMAN F, MERT A C, ÖZTÜRK E, et al. A hardware accelerator for polynomial multiplication operation of CRYSTALS-Kyber PQC scheme[C]//Proceedings of 2021 Design, Automation & Test in Europe Conference & Exhibition (DATE). Piscataway: IEEE Press, 2021: 1020-1025.
- [20] KUO P C, CHEN Y W, HSU Y C, et al. High performance post-quantum key exchange on FPGAs[J]. Journal of Information Science & Engineering, 2021, 37(5): 1211-1229.
- [21] 李斌, 陈晓杰, 冯峰, 等. 后量子密码 CRYSTALS-KYBER 的 FPGA 多路并行优化实现[J]. 通信学报, 2022, 43(2): 196-207.
LI B, CHEN X J, FENG F, et al. FPGA multi-unit parallel optimization and implementation of post-quantum cryptography Crystals-Kyber[J]. Journal on Communications, 2022, 43(2): 196-207. (in Chinese)
- [22] XING Y F, LI S G. An efficient implementation of the NewHope key exchange on FPGAs[J]. IEEE Transactions on

Circuits and Systems I: Regular Papers, 2020, 67(3): 866-878.

- [23] BISHEH-NIASAR M, AZARDERAKH-SH R, MOZAFARI-KERMANI M. Instruction-set accelerated implementation of Crystals-Kyber[J]. IEEE Transactions on Circuits and Systems I: Regular Papers, 2021, 68(11): 4648-4659.

作者简介



吕顺森 男, 1999 年出生, 河南信阳人. 郑州大学硕士生. 主要研究方向为后量子密码、高性能计算.
E-mail: l_sh_sen@163.com



李斌 男, 1986 年出生, 河南郑州人. 博士. 郑州大学讲师. 主要研究方向为信息安全、可重构计算.
E-mail: iebinli@zzu.edu.cn



翟嘉琪 男, 1998 年出生, 河南郑州人. 郑州大学硕士生. 主要研究方向为高性能计算、体系结构.
E-mail: zhajiaqi777@163.com



李松岐 男, 1998 年出生, 河南洛阳人. 郑州大学硕士生. 主要研究方向为高性能计算、体系结构.
E-mail: zzuqiezz@163.com



周清雷 男, 1962 年出生, 河南新乡人. 博士. 郑州大学教授. 主要研究方向为信息安全、自动机理论和计算复杂性理论.
E-mail: ieqlzhou@zzu.edu.cn