Gadget-based Masking of Streamlined NTRU Prime Decapsulation in Hardware

Georg Land¹, Adrian Marotzke^{2,3}, Jan Richter-Brockmann¹ and Tim Günevsu^{1,4}

Abstract. Streamlined NTRU Prime is a lattice-based Key Encapsulation Mechanism (KEM) that is, together with X25519, the default algorithm in OpenSSH 9. Based on lattice assumptions, it is assumed to be secure also against attackers with access to large-scale quantum computers. While Post-Quantum Cryptography (PQC) schemes have been subject to extensive research in recent years, challenges remain with respect to protection mechanisms against attackers that have additional side-channel information, such as the power consumption of a device processing secret data. As a countermeasure to such attacks, masking has been shown to be a promising and effective approach. For public-key schemes, including any recent PQC schemes, usually, a mixture of Boolean and arithmetic techniques is applied on an algorithmic level. Our generic hardware implementation of Streamlined NTRU Prime decapsulation, however, follows an idea that until now was assumed to be solely applicable efficiently to symmetric cryptography: gadget-based masking. The hardware design is transformed into a secure implementation by replacing each gate with a composable secure gadget that operates on uniform random shares of secret values. In our work, we show the feasibility of applying this approach also to PQC schemes and present the first Public-Key Cryptography (PKC) – pre- and post-quantum – implementation masked with the gadget-based approach considering several trade-offs and design choices. By the nature of gadget-based masking, the implementation can be instantiated at arbitrary masking order. We synthesize our implementation both for Artix-7 Field-Programmable Gate Arrays (FPGAs) and 45 nm Application-Specific Integrated Circuits (ASICs), yielding practically feasible results regarding the area, randomness requirement, and latency. We verify the side-channel security of our implementation using formal verification on the one hand, and practically using Test Vector Leakage Assessment (TVLA) on the other. Finally, we also analyze the applicability of our concept to Kyber and Dilithium, which will be standardized by the National Institute of Standards and Technology (NIST).

Keywords: PQC, Masking, FPGA, ASIC, Streamlined NTRU Prime, Higher-order Masking, Gadget-based Masking

1 Introduction

Wide deployment of Post-Quantum Cryptography (PQC) algorithms in devices and applications is indispensable, even though there is no guarantee that the advent of large-scale quantum computers will happen at all. For many application scenarios, security-critical implementations must also provide security against physical attacks.

¹ Ruhr University Bochum, Horst Görtz Institute for IT Security, Bochum, Germany mail@georg.land,{jan.richter-brockmann,tim.gueneysu}@rub.de

² Hamburg University of Technology, Hamburg, Germany adrian.marotzke@tuhh.de
³ NXP Semiconductors, Hamburg, Germany adrian.marotzke@nxp.com
⁴ DFKI GmbH, Cyber-Physical Systems, Bremen, Germany

Embedded devices in particular require considering adversaries who can measure timing, the power consumption or electromagnetic (EM) emanation of a device processing secret data. In this context, many practical attacks have been shown in the past on "classical", but also PQC schemes. For instance, several attacks have been published attacking Kyber [XPR+21,SPH22,HHP+21,HPP21], Saber [NDGJ21], Falcon [KA21], NTRU [AR21], or even generic on lattice-based constructions [RRCB20]. Notable here are attacks targeting side-channel protected implementations, such as the recent attacks on a fifth-order masked Kyber implementation [DNG22], a third-order masked Saber implementation [NWDP22] and a first-order masked Saber implementation [NDJ21].

Specifically for Streamlined NTRU Prime, two attacks have been proposed. First, Xu et al. show single-trace attacks on fixed weight sampling as used in Streamlined NTRU Prime and NTRU key generation as well as Dilithium signing [KAA21]. Furthermore, Ravi et al. present a method to recover the Streamlined NTRU Prime secret key with a side-channel assisted chosen-ciphertext attack [FBR $^+$ 22]. They demonstrate the capability of a full key recovery with just 3 005 traces for the smallest parameter set and with 4 688 traces for larger parameter sets.

Contrary, dedicated countermeasures aiming at decoupling the connection between secret data and power consumption have been proposed in the past decades. The primary technique for that purpose is *masking*, which splits secret values into multiple uniform random shares. In this context, research has recently focused on masked PQC implementations in *software*, mainly for Kyber and Saber. A recent work presents a masked implementation of NTRU for embedded software [CGTZ23]. In contrast, there are far fewer works on hardening hardware implementations, again focusing on Kyber [FBR⁺22, KNAH22] and Saber [AMD⁺21, FBR⁺22]. To the best of our knowledge, no PQC schemes other than Kyber and Saber have a fully masked hardware implementation published, where both implementations target Field-Programmable Gate Arrays (FPGAs) and first-order security only.

To date, secure implementations for the Streamlined NTRU Prime scheme have not yet been proposed, neither for software nor for hardware platforms, despite the fact Streamlined NTRU Prime is already the default choice for the widely used OpenSSH suite, starting from version 9.0. In our work, we aim to close this gap by presenting the first masked hardware implementation of Streamlined NTRU Prime decapsulation for hardware devices, synthesizable both for FPGA and Application-Specific Integrated Circuit (ASIC) and aiming at use-cases where OpenSSH-supported connections are established with external hardware devices that are potentially under exposure to physical attackers.

Moreover, our implementation is masked on the *gate level* rather than the algorithmic level, which has the advantage of being easily configurable to provide protection for any arbitrary order. To the best of our knowledge, this is the first gadget-based masked implementation of a Public-Key Cryptography (PKC) decapsulation or decryption, both for the pre- and post-quantum settings, and the first masked ASIC implementation of any PQC scheme. Thus, our work provides a milestone for the practicality of gadget-based masking, extending its potential application space to PKC and PQC in particular.

Contribution We provide the following contributions:

- This work presents for the first time gadget-based masked implementations of any PKC scheme.
- Our approach can be generalized to an arbitrary-order masked hardware implementation of any PQC scheme for which the masking degree can be adjusted easily.
- We implement our design both on a Xilinx Artix-7 FPGA, and as an ASIC using the 45 nm Nangate open cell library¹.

¹Available at https://si2.org/open-cell-library/

- Compared to other existing fully masked PQC FPGA implementations, our implementation has similar (in the case of Saber) or significantly lower (in the case of Kyber) resource requirements for first-order security.
- We present the first arbitrary-order masked SHA-2 hardware implementation in literature.
- The side-channel resistance of our implementation is formally verified using VER-ICA [RFSG22] and practical measurements.
- Our source code is publicly available at https://github.com/AdrianMarotzke/Masked-SNTRUP.

2 Preliminaries

In this section, we briefly introduce the notations used throughout this work. Afterward, we recap masking and important composability notions. Eventually, we describe Streamlined NTRU Prime and particularly the decapsulation.

2.1 Notation

Throughout this work, we denote $\mathcal{R}_q = \mathbb{Z}_q[x]/(x^p-x-1)$, and $\mathcal{R}_3 = \mathbb{Z}_3[x]/(x^p-x-1)$ with p,q being primes. Furthermore, we write x[i:j] for bit vectors of length |i-j|+1 and also allow multiple dimensions for this, e.g., x[i:j,k:l] is a vector of |i-j|+1 bit vectors each of length |k-l|+1. For masking, we use d as the masking degree, i.e., the number of probes an attacker has access to. It follows that we split secrets into d+1 shares, referring to a single share as $x^{(i)}$ with $0 \le i \le d$. Moreover, we denote $x^{(0:d)}$ as a masked variable. At any occurrence of Boolean operations that involve masked variables, we assume to perform this securely, e.g., by means of a secure gadget. Finally, we stress that $\overline{x^{(0:d)}}$ denotes inverting the secret value by inverting one share rather than inverting each share (which would not invert the secret value for odd d).

2.2 Masking

Masking is an approach based on Shamir's secret sharing. It has been proven as an effective countermeasure against power or EM side-channel attacks by splitting secret values into uniform random shares. In our work, we employ only Boolean masking where a secret value x is split into d+1 shares $x^{(i)}$, such that $x=\bigoplus_{i=0}^d x^{(i)}$. While functions that are linear or affine in the masking domain can be applied trivially to each share individually, we use specialized methods to secure non-linear functions like AND or OR operations.

The core concept of gadget-based masking is to replace individual hardware gates with secure versions, so-called gadgets. These secure hardware gates are designed to not leak their inputs and outputs via power consumption or EM emanation. Early versions of these secure gates aimed at ensuring that the power consumption remained constant, regardless of the input or output. An example of such an approach is Wave Dynamic Differential Logic (WDDL) [TV04], which is a type of Dual-Rail with Precharge (DRP) logic. These logic gates use differential inputs and outputs and have a pre-charge phase which ensures that transistors switch at every clock cycle, even if the inputs do not change. An overview of other DRP logic styles can be found in [DGBN09]. However, many of the DRP gates were successfully attacked over the years [SGD+09, MKEP11, PKZM07, DGBN09]. These attacks were possible due to effects such as unbalanced routing of the differential signals or glitches in the circuit.

In order to properly evaluate and formally verify the resistance against side-channel attacks of such special gates, a range of different attacker models have been proposed in

the past. In 2003, Ishai, Sahai, and Wagner [ISW03] introduced the d-probing model, which is still frequently used as an appropriate abstraction. However, this model neither includes glitches nor transitions or couplings and thus has been extended to the *robust* d-probing model incorporating these phenomena [BGI⁺18, FGP⁺18].

Nevertheless, the robust d-probing model is insufficient to analyze the composability of gadgets (nowadays also called gate-level masking). Hence, Barthe et al. introduced Non-Interference (NI) as the first composability notion in 2015 [BBD+15]. Although NI limits the leakage between shared intermediate results, it does not guarantee probing security of composed circuits. Therefore, Barthe et al. presented the notion of Strong Non-Interference (SNI) [BBD+16], which ensures composability of gadgets. Eventually, Cassiers and Standaert proposed Probe-Isolating Non-Interference (PINI) [CS20] reducing the overhead introduced by SNI gadgets. PINI ensures that all shared AND gadgets are composable and XOR as well as NOT operations can be performed share-wise without refreshing.

Bringing this concept to concrete instantiations of SecAND gadgets in hardware, Cassiers et al. proposed Hardware Private Circuit (HPC) [CGLS21]. HPC allows instantiating an arbitrary-order masked SecAND gadget with two clock cycles latency for one input and one clock cycle for the other input denoted as HPC1. Moreover, they optimized this gadget for less randomness demand denoted as HPC2 gadget. Following this, Knichel et al. proposed Generic Hardware Private Circuits (GHPCs) to build more complex PINI gadgets [KSM22]. Finally, in a recent work, Knichel and Moradi presented HPC3 achieving lower latency by using more fresh randomness [KM22].

2.3 Streamlined NTRU Prime

Streamlined NTRU Prime is a lattice-based Key Encapsulation Mechanism (KEM) that is resistant against both classical and quantum adversaries [BCLv17, BBC $^+$ 20]. It has been designed carefully using structured lattices while firmly avoiding potentially exploitable attack surfaces. In particular, it eliminates decryption failures and employs large Galois groups instead of cyclotomics.

Streamlined NTRU Prime defines Short as the set of polynomials in \mathcal{R}_q with exactly w non-zero coefficients from $\{-1,1\}$. Furthermore, we also use the notation of an underline indicating that the respective value is encoded.

As a KEM, it uses the Fujisaki-Okamoto transform to achieve indistinguishability under chosen-ciphertext attacks (IND-CCA) and builds upon a public-key encryption scheme that fulfills one-wayness against passive attacks. In the following, we describe the three procedures of the KEM: key generation, encapsulation, and decapsulation.

Key Generation. First, a uniform random polynomial g in \mathcal{R}_3 is generated. This step is repeated until g is invertible in \mathcal{R}_3 . Then, the inverse polynomial of g is computed. Furthermore, f is sampled to be a polynomial from Short. The secret key consists of f and g^{-1} as well as a random bit string ρ which is used for implicit rejection during decapsulation. Finally, the public key is computed as h = g/(3f) where $h \in \mathcal{R}_q$.

Encapsulation. The first step is to sample a uniformly random polynomial r from Short, which is then multiplied with the public key polynomial h. In the resulting polynomial, each coefficient is rounded to the nearest multiple of three. The output of this operation is denoted as the polynomial c. Subsequently, the encoded r and the encoded public key are hashed to create the ciphertext confirmation hash. The confirmation hash together with the encoded c is the ciphertext. The session key is computed by hashing the encoded r and the ciphertext.

Algorithm 1 Streamlined NTRU Prime Decapsulation [BBC+20]

```
Require: ciphertext C = (c, \gamma), secret key (k = \mathsf{Encode}(f, g^{-1}), K = \mathsf{Encode}(h), \rho,
      hash_4(K)
 1: c \in \mathcal{R}_3 := \mathsf{Decode}(\underline{c})
 2: (f, v) \in \mathcal{R}_3 \times \mathcal{R}_3 := \mathsf{Decode}(\underline{k})
 3:\ h\in\mathcal{R}_q:=\mathsf{Decode}(\underline{K})
 4: e \in \mathcal{R}_3 := ((3fc) \in \mathcal{R}_q) \mod 3
 5: r' \in \mathcal{R}_3 := ev
 6: if r' does NOT have weight w then
           r' := (1, 1, \dots, 1, 0, 0, \dots, 0)
                                                                              \triangleright The first w elements are 1, the rest 0
 8: c' \in \mathcal{R}_q := \mathsf{Round}(hr')
                                                            \triangleright re-encrypt with h, r', compute new ciphertext c'
 9: \underline{c}' := \mathsf{Encode}(c')
10: \underline{r}' := \mathsf{Encode}(r')
11: \gamma' := \mathsf{hash}_2(\mathsf{hash}_3(\underline{r}'), \mathsf{hash}_4(\underline{K}))
                                                                 > re-compute the ciphertext confirmation hash
12: C' = (\underline{c'}, \gamma')
13: if C' = C then
           return hash_1(hash_3(\underline{r}'), C)
14:
15: else
           return hash_0(hash_3(\rho), C)
16:
```

Decapsulation. The decapsulation is shown in Algorithm 1 in detail. The basic idea is to remove the denominator of the public key from the ciphertext by multiplying 3f in \mathcal{R}_q . The subsequent application of modulo 3 to each coefficient removes the rounding error which is succeeded by the multiplication with $1/g \in \mathcal{R}_3$ to also remove the numerator of the public key and to obtain the plaintext. This plaintext is checked to be in the correct space Short. Furthermore, to ensure that no chosen-ciphertext attack is carried out, the obtained plaintext is re-encrypted and the result is compared to the original ciphertext. If everything matches, the correct session key is reconstructed, else an implicit rejection is performed by using ρ . Note that this final rejection step is strictly required to be performed in constant time.

3 Conceptual Considerations

To implement the decapsulation as shown in Algorithm 1, we essentially need six major modules:

- 1. Polynomial multiplication with operands in $(\mathcal{R}_q, \mathcal{R}_3)$ and return values in \mathcal{R}_q ,
- 2. Polynomial multiplication with operands in $(\mathcal{R}_3, \mathcal{R}_3)$ and result in \mathcal{R}_3 ,
- 3. Reduction component modulo 3,
- 4. Weight check component,
- 5. Rounding module, and
- 6. SHA-512.

Standard Approach. Usually, to mask polynomial multiplication modules, additive masking would be applied, with either multiple polynomial multipliers being instantiated in parallel, or one polynomial multiplier being instantiated that processes the shares consecutively. Moreover, two of the three multiplications have one public and one secret input, which can be realized very efficiently by applying additive masking as it only requires d+1 polynomial multiplications and no re-sharing. The other multiplication, however,

has two secret input polynomials. In order to perform a secure polynomial multiplication in the additive domain, $\frac{d^2+d}{2}$ fresh random polynomials need to be sampled. Additionally, $2(d^2+d)$ polynomial additions and d^2+d polynomial multiplications must be performed.

In contrast, masking the reduction, weight check, and rounding is non-trivial in the arithmetic domain and would be solved in the Boolean domain. Finally, SHA-512 uses 64 bit additions, which is efficient in additive domain and feasible but less efficient in Boolean domain, as well as non-linear Boolean operations that strictly require Boolean masking.

In summary, this traditional approach is expected to yield a relatively efficient implementation at the cost of converting between additive and Boolean masking domain multiple times. Moreover, this type of implementation is often very specific in terms of masking degree, i.e., not being parametrizable. Besides, the wide variety of applied techniques produces a larger attack surface, as shown in recent attacks on masking conversions [NWDP22].

Applicability of Gadget-based Masking. To overcome these downsides, we follow a recent line of research from the field of masking symmetric cryptographic schemes: gadget-based masking. For schemes in symmetric cryptography, we usually find a Boolean description that enables masking them at the gate level. This differs for public-key and post-quantum cryptography as these schemes typically employ arithmetic operations on number-theoretic structures such as multiplications in polynomial fields. Polynomial multiplications, however, consist of modular multiplications and additions in some finite number field. While the modular additions can be masked easily in Boolean domain through a secure adder, the modular multiplications are vastly more complex and are deemed infeasible to be masked in the Boolean domain.

However, for Streamlined NTRU Prime, we observe that the three polynomial multiplications each have at least one factor in \mathcal{R}_3 . Consequently, if we employ schoolbook multiplication, the underlying coefficient multiplication-accumulation has an input from \mathbb{Z}_q being multiplied either 1, 0, or -1 and then accumulated to another value in \mathbb{Z}_q . We immediately observe that no complex modular multiplication must be carried out in this case. Instead, we can securely multiplex between the input coefficient from \mathbb{Z}_q , its precomputed additive inverse, and zero. The result is added securely to the accumulation value. As indicated before, all other operations are already feasible in Boolean domain, enabling the first fully Boolean masked implementation of a public key and post-quantum secure scheme.

In the following, we describe our design considerations for each module in Boolean domain. Note that in contrast to conventional hardware development, where it is desirable to have as many NAND gates as possible as they are the smallest gates, the design goal in our case is to have as few as possible SecAND gadgets, as they require fresh randomness. Throughout our design, we use the HPC2 SecAND gadget [CGLS21].

3.1 Polynomial Multiplication

Polynomial multiplications are the most expensive operations in decapsulation. Thus, research usually focuses on improving their performance [Mar20, PMT+22, CHK+21, ACC+21]. Instead, we focus on achieving a *secure* implementation. During decapsulation, two types of multiplications are required: 1. Multiplication in \mathcal{R}_q with one operand from \mathcal{R}_3 (Lines 4 and 8 in Algorithm 1) and 2. Multiplication in \mathcal{R}_3 (Line 5 in Algorithm 1).

3.1.1 Multiplication in \mathcal{R}_q

We observe that if we employ a standard schoolbook multiplication approach for both occasions of this multiplication, no coefficient multiplier is necessary. Instead, we use a secure adder and a secure three-way multiplexer. It is important to note that for both

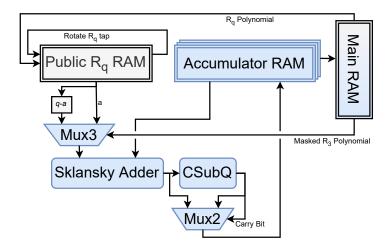


Figure 1. Architecture of the \mathcal{R}_q polynomial multiplier. Blue modules operate on masked shares.

multiplications in \mathcal{R}_q , the input polynomial from \mathcal{R}_q is public, while the other factor from \mathcal{R}_3 is secret. Thus, the idea is to compute the additive inverse of the input coefficient from \mathcal{R}_q , which is unmasked. Then, we securely multiplex with the masked select signal between both values and zero, and finally accumulate the result securely to the (intermediate) result coefficient. The architecture is shown in Figure 1.

Secure Multiplexing. Furthermore, we need a secure three-way multiplexer. The three *public* input signals are $z=0, a_p=a, a_n=q-a\in\mathbb{Z}_q$. However, here we view them as Boolean values in \mathbb{F}_2^{13} . The secret select signal is $(f[1], f[0]) \in \{(0,0), (0,1), (1,1)\}$. We perform two consecutive secure 2-input multiplexing operations:

$$x[0]^{(0:d)} = a_p \wedge f[1]^{(0:d)} \oplus a_n \wedge \overline{f[1]^{(0:d)}} = a_p \wedge f[1]^{(0:d)} \oplus a_n \wedge (f[1]^{(0:d)} \oplus 1)$$

$$= a_p \wedge f[1]^{(0:d)} \oplus a_n \wedge f[1]^{(0:d)} \oplus a_n$$

$$= ((a_p \oplus a_n) \wedge f[1]^{(0:d)}) \oplus a_n \qquad (1)$$

$$x[1]^{(0:d)} = x[0]^{(0:d)} \wedge f[0] \oplus z \wedge \overline{f[0]^{(0:d)}} = ((x[0]^{(0:d)} \oplus z) \wedge f[0]^{(0:d)}) \oplus z$$

$$= x[0]^{(0:d)} \wedge f[0]^{(0:d)} \qquad (2)$$

Note that the public inputs can be set as first shares and all other shares are just zeros. This is the reason why we can simply omit z in Equation 2. The SecAND gadget generates a uniformly random output also for the case that $(f_1, f_0) = (0, 0)$.

Secure Addition. Parallel prefix adders can achieve efficient addition in hardware. These concepts also have been adapted to the Boolean masked domain first in [SMG15]. This was followed by a broader examination of more recent techniques like threshold implementation and gadget-based masking [BG22], which we deploy for our work.

3.1.2 Multiplication in \mathcal{R}_3

For mulitplications in \mathcal{R}_3 only nine possible input combinations with three output combinations exist. Thus, we develop a direct Boolean masking utilizing the fact that the single inputs have a limited range. Multiplying two signed two-bit coefficients e[1:0] = e[0] - 2e[1] and v[1:0] = v[0] - 2v[1] to a signed two-bit value r[1:0] = r[0] - 2r[1] can be done as

follows:

$$r[0]^{(0:d)} = e[0]^{(0:d)} \wedge v[0]^{(0:d)}$$
(3)

$$r[1]^{(0:d)} = e[0]^{(0:d)} \wedge v[0]^{(0:d)} \wedge (e[1]^{(0:d)} \oplus v[1]^{(0:d)})$$

$$\tag{4}$$

Then, we add $r[1:0]^{(0:d)}$ to the accumulation value $a[1:0]^{(0:d)}$ and map the result back to the signed $a'[1:0]^{(0:d)} \in \{-1,0,1\}$ which can be done with the following formulas that take into account that only $00_2,01_2,11_2$ are valid inputs:

$$a'[0]^{(0:d)} = \left(r[0]^{(0:d)} \oplus a[0]^{(0:d)}\right) \vee \left(r[0]^{(0:d)} \wedge \left(\overline{r[1]^{(0:d)} \oplus a[1]^{(0:d)}}\right)\right) \tag{5}$$

$$a'[1]^{(0:d)} = \left(r[1]^{(0:d)} \wedge \overline{a[0]^{(0:d)}}\right) \oplus \left(\overline{r[1]^{(0:d)}} \wedge \left(r[0]^{(0:d)} \oplus a[1]^{(0:d)}\right)\right) \tag{6}$$

3.1.3 Schoolbook Polynomial Multiplication

Generally, there are three approaches for this: Either we rotate one of the input polynomials or the output polynomial. For our two "big" multiplications in \mathcal{R}_q , we have a small secret input represented by 2(d+1) bits, a big public input represented by $\lceil \log_2 q \rceil$ bits, and a big secret output represented by $(d+1)\lceil \log_2 q \rceil$ bits. Since shifting many data is expensive in terms of routing, Flip-Flop (FF) demand, and dynamic power consumption, the natural choice is to rotate either of the input polynomials.

3.1.4 Polynomial Reduction modulo $x^p - x - 1$

For the schoolbook multiplication, we can directly perform the polynomial reduction. We observe that $x^p \equiv x+1 \mod x^p-x-1$, which indicates that the uppermost coefficient (x^p) during rotation must be additionally added to the before lowermost coefficient. As we indicated before, we want to rotate either of the input polynomials. Applying this strategy to the \mathcal{R}_3 polynomial would increase the coefficient range to [-2,2] due to the extra addition during polynomial reduction. We would require a 5-way multiplexer instead of a 3-way multiplexer, increasing both area and randomness demand. Thus, we choose to rotate the public \mathcal{R}_q input polynomial and perform the polynomial reduction in the same domain.

3.2 Modular Reductions

For Streamlined NTRU Prime decapsulation, we require two different modular reductions.

Reduction Modulo q. This reduction is only applied for the accumulation within the \mathcal{R}_q polynomial multiplications. We decided to use the non-negative modular representation in the interval [0,q) only since we would need to check both for underflows and overflows in the centered representation. Therefore, the value to reduce only grows by a maximum of one bit and can only provoke an overflow. Thus, a conditional subtraction by q suffices, which we perform as follows.

We subtract q from all accumulation results and obtain the carry bit from that subtraction. If this is 1, we know an underflow occurred. Thus, we can use the carry bit to securely multiplex between the original accumulation value and the subtracted value. This keeps all intermediate values in the minimal interval [0, q).

Reduction Modulo 3. For the modulo 3 reduction, we have given an input from \mathbb{Z}_q and want to reduce it to $\{-1,0,1\}$. We start with an unsigned 13-bit number z[12:0] and repeatedly exploit the relation $2 \equiv -1 \mod 3$. Note that all operations are carried out in

masked domain, but we omit the masking notation when dealing with arithmetic modulo 3.

$$z[12:0] \equiv 2z[12:1] + z[0] \equiv -z[12:1] + z[0] \mod 3$$

$$\equiv -2z[12:2] - z[1] + z[0] \equiv z[12:2] - z[1] + z[0] \mod 3$$

$$\equiv \sum_{i=0}^{6} z[2i] - \sum_{i=0}^{5} z[2i+1] \mod 3$$
(7)

The result of this computation ranges from -6 to 7 and is represented by a signed 4-bit integer $y[3:0] = -2^3y[3] + y[2:0]$. We again exploit the above relation:

$$-2^{3}y[3] + y[2:0] \equiv y[3] + y[2:0] \equiv y[3] + 2y[2:1] + y[0] \equiv y[3] - y[2:1] + y[0] \mod 3$$

$$\equiv y[3] - 2y[2] - y[1] + y[0] \equiv y[3] + y[2] - y[1] + y[0] \mod 3$$
 (8)

This results in a value ranging from -1 to 3, represented by a signed 3-bit integer $x[2:0] = -4x[2] + x[1:0] \equiv x[1:0] - x[2] \mod 3$. This value can already be mapped to a value $w[1:0] \in \{-1,0,1\}$ efficiently:

$$w[0]^{(0:d)} = x[0]^{(0:d)} \oplus x[1]^{(0:d)} \oplus x[2]^{(0:d)}$$
(9)

$$w[1]^{(0:d)} = \left(x[0]^{(0:d)} \wedge x[1]^{(0:d)}\right) \oplus x[1]^{(0:d)} \oplus x[2]^{(0:d)}$$
(10)

One additional point to consider is that this modulo 3 calculation assumes an unsigned 13-bit number. However, in the NTRU Prime specification, the modulo 3 operation is used on signed 13-bit numbers, in the interval [-q/2, q/2] [BBC⁺20]. This means that numbers in the interval [q/2, q) must be treated slightly differently, as these were originally negative. However, the solution is simple: since q = 4591, and $4591 = 1 \mod 3$, we simply have to add 1 to the final result if the original number was in the interval [q/2, q). This addition can be in a similar way to the multiplication in \mathcal{R}_3 (see Section 3.1.2).

3.3 Weight Check

Let $r'[0:1,0:p-1]^{(0:d)}$ be an array of p shared two-bit numbers. Valid values are (0,0), (0,1), (1,1) if the signed representation is used, and (0,0), (0,1), (1,0) for the unsigned representation. We wish to check if exactly w array elements are non-zero. Thus, the basic idea depends on the chosen representation.

Signed. We accumulate all $r'[0,:]^{(0:d)}$ values together, with a secure $\lceil \log_2 w \rceil$ -bit adder.

Unsigned. We compute $r'[0,:]^{(0:d)} \vee r'[1,:]^{(0:d)}$ and accumulate the resulting shared bit vector with a $\lceil \log_2 w \rceil$ -bit adder.

It follows that the signed representation demands fewer non-linear Boolean operations. For the secure adder, the same adder as used for the polynomial multiplications is applied.

Following this, we then bit-wise XOR the shared adder output with the public target weight w, and then OR all bits of the result together to a single shared result bit.

The overwriting of r' can be performed with a secure 2-way multiplexer deciding between the secret r' and the fixed public vector (1, 1, ..., 1, 0, 0, ..., 0).

3.4 Rounding

For rounding, we first perform a reduction of the coefficient modulo 3 and then subtract the result from the original coefficient. As a result of the modulo operation, we obtain two masked bits $a[1:0]^{(0:d)} \in \{(0,0),(0,1),(1,1)\}$. With this, we want to

- 1. add 1 for $a[1:0]^{(0:d)} = (1,1)$
- 2. add q-1, which is analogue to subtracting 1, for $a[1:0]^{(0:d)}=(0,1)$, and
- 3. add zero for $a[1:0]^{(0:d)} = (0,0)$

One way to achieve this is by multiplexing securely between q-1, 1 and zero depending on $a[1:0]^{(0:d)}$, which in turn would include more non-linearity. To avoid this, we can construct the value $a[0]^{(0:d)} \cdot (q-1) - a[1]^{(0:d)} \cdot q$ and add that to the initial coefficient. In other words, this value consists of $a[0]^{(0:d)}$ in all binary positions where q-1 is 1, except the least significant bit, where it consists of $a[0]^{(0:d)} \oplus a[1]^{(0:d)}$. For the addition, we can re-use the addition-reduction procedure as used for polynomial multiplication.

3.5 SHA-512

SHA-512 employs a Merkle-Damgård construction processing a 512 bit state divided into eight 64 bit words A, B, C, D, E, F, G, H. In order to update the state, SHA-512 implements seven adders (modulo 2^{64}), the two functions Σ_0 and Σ_1 , and the functions SHA-Ch and SHA-Ma. The former two functions Σ_0 and Σ_1 consist of simple shift operations by three different values for each function processing A and E, respectively. The outputs of the shifts are added together by XOR operations. SHA-Ch and SHA-Ma are non-linear functions processing E, F, G and A, B, C, respectively.

For our masked hardware implementation, we protect the seven adders by applying the concept of the masked adder introduced in Section 3.1. We instantiate a complete 64-bit adder to realize the correct addition. Masking Σ_0 and Σ_1 can be accomplished in a straightforward way since the shift operations do not introduce additional implementation overhead in hardware and all XOR gates can simply be replaced by secure XOR gadgets.

Finally, SHA-Ch and SHA-Ma are bit-wise operations that can be implemented in parallel to match the width of the adder to be used. Hence, we can modify the formulas for both to reduce the number of non-linear gates in order to minimize the amount of required randomness and the area overhead leading to

$$\begin{aligned} \mathsf{SHA-Ch}(E,F,G) &= (E \wedge F) \oplus (\overline{E} \wedge G) = (E \wedge F) \oplus ((E \oplus 1) \wedge G) \\ &= (E \wedge (F \oplus G)) \oplus G \\ \mathsf{SHA-Ma}(A,B,C) &= (A \wedge B) \oplus (A \wedge C) \oplus (B \wedge C) \\ &= (A \wedge (B \oplus C)) \oplus (B \wedge C). \end{aligned} \tag{11}$$

3.6 Encoding, Decoding & Comparison

Streamlined NTRU Prime defines multiple en- and decoding algorithms for transforming polynomials in \mathcal{R}_3 and \mathcal{R}_q to and from byte arrays [BBC⁺20]. Decoding the ciphertext and public key can be done unmasked as they are both public. We use the decoder described in [PMT⁺22]. For decoding the secret polynomials f and g^{-1} , we also use the decoder from [PMT⁺22], and apply masking afterwards. However, we need to securely encode r' into a byte array to compute the confirmation hash and session key. For this, we apply masking to the \mathcal{R}_3 encoder from [PMT⁺22]. This is straightforward as the encoder only consists of a shift register and a 2-bit adder.

In the original algorithm specification, the recomputed ciphertext polynomial c' is encoded (line 9 in Algorithm 1) before the ciphertext comparison (line 13), using an \mathcal{R}_q encoder. However, the \mathcal{R}_q encoder requires a 16-bit multiplication which would be prohibitively expensive to implement securely. We instead compare the ciphertext polynomial coefficients directly, after which we compare the confirmation hashes. This prevents us from implementing the masked \mathcal{R}_q encoder. The masked ciphertext comparison is straightforward: We do a bit-wise secure XOR of the two ciphertext coefficients and then repeatedly OR the output together.

4 Implementation

After introducing the theoretical background of masking all required operations, we now discuss the implementations of each building block.

Add13 and Add64. In their work [BG22], Bache and Güneysu compare the Brent-Kung, Kogge-Stone, and Sklansky adder architectures in the context of Boolean masking. For gadget-based masking, the Sklansky adder turns out to be the optimal choice, having the same low latency as Kogge-Stone but less randomness demand while having a lower latency than Brent-Kung at the cost of slightly more randomness.

The 13-bit Sklansky adder with carry-out deployed in our implementation is shown in Figure 2a. For input bits $a[i]^{(0:d)}$, $b[i]^{(0:d)}$ where $i \in \{0, ..., 12\}$, we compute in each circle:

$$g[i]^{(0:d)} = a[i]^{(0:d)} \wedge b[i]^{(0:d)}$$
(13)

$$p[i]^{(0:d)} = a[i]^{(0:d)} \oplus b[i]^{(0:d)}$$
(14)

Note that the dotted circle indicates that the input is all zero, and requires no computation. It is needed to compute the carry-out, as we have a 14-bit output, and ensures that the lower layers operate on the correct inputs. Each square node has four inputs, the two "left" inputs $g_l^{(0:d)}, p_l^{(0:d)}$ and the two "right" inputs $g_r^{(0:d)}, p_r^{(0:d)}$, and computes the following outputs:

$$g^{(0:d)} = g_l^{(0:d)} \oplus \left(p_l^{(0:d)} \wedge g_r^{(0:d)} \right)$$
 (15)

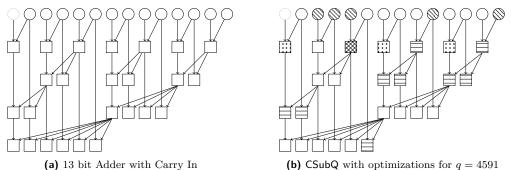
$$p^{(0:d)} = p_l^{(0:d)} \wedge p_r^{(0:d)} \tag{16}$$

Finally, note that all leaf nodes do not need to compute $p^{(0:d)}$, as only the final $g^{(0:d)}$ values are needed.

The 64-bit adder works equivalently, though with a total of six levels. In this case, we do not need a carry-in or carry-out.

CSubQ. For the conditional subtraction with q, we take a similar approach. We instantiate another Sklansky adder with one public operand fixed to the two's complement of q. Then, after each addition (let us denote the result here as $x^{(0:d)}$), we perform this subtraction by q and obtain $(q-x)^{(0:d)}$ as well as the shared carry-out bit. Using this, we multiplex securely between $x^{(0:d)}$ and $(q-x)^{(0:d)}$, selecting the former if the carry-out is one (indicating an underflow has occurred) and else the latter one.

The fixed input already enables vast optimizations by the synthesizer. Further improvements could be made by optimizing the adder architecture itself for a fixed operand. Since we know the positions of the zeros, we could simplify our adder as depicted in Figure 2b.



(a) 13 bit Adder with Carry In (b) CSubQ with c Figure 2. Sklansky Adder Constructions

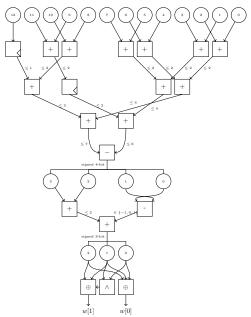


Figure 3. Mod3 module

However, note that we did not implement these optimizations and have left them for future work.

The computation of all p values below the first row is the same as before. However, we can completely omit computing the first row of p, q as described in Equation 13 and Equation 14. Instead, we know, given an input $a[12:0]^{(0:d)}$, for each circle in Figure 2b that

$$g[i]^{(0:d)} = \begin{cases} a[i]^{(0:d)} & \text{if } (-q)[i] = 1\\ 0 & \text{else} \end{cases}$$

$$p[i]^{(0:d)} = \begin{cases} \overline{a[i]^{(0:d)}} & \text{if } (-q)[i] = 1\\ a[i]^{(0:d)} & \text{else} \end{cases}$$
(17)

$$p[i]^{(0:d)} = \begin{cases} \overline{a[i]^{(0:d)}} & \text{if } (-q)[i] = 1\\ a[i]^{(0:d)} & \text{else} \end{cases}$$
 (18)

In Figure 2b, the circles filled with the diagonal line pattern indicate that the fixed input bit of the two's complement of q is one. For the squares, we have four different cases now:

Non-filled Computed as before.

$$\mathbf{Grid} \hspace{1cm} g^{(0:d)} = g_l^{(0:d)} \oplus \left(p_l^{(0:d)} \wedge g_r^{(0:d)} \right) = g_l^{(0:d)} \oplus \left(p_l^{(0:d)} \wedge 0 \right) = g_l^{(0:d)}$$

$$\textbf{Dotted} \qquad \qquad g^{(0:d)} = g_l^{(0:d)} \oplus \left(p_l^{(0:d)} \wedge g_r^{(0:d)} \right) = 0 \oplus \left(p_l^{(0:d)} \wedge 0 \right) = 0$$

$$\text{Horizontal lines} \qquad g^{(0:d)} = g_l^{(0:d)} \oplus \left(p_l^{(0:d)} \wedge g_r^{(0:d)} \right) = 0 \oplus \left(p_l^{(0:d)} \wedge g_r^{(0:d)} \right) = p_l^{(0:d)} \wedge g_r^{(0:d)}$$

Mod3 and Mul3. The architecture to compute Mod3 is depicted in Figure 3. For the secure additions and subtractions, we employ simple ripple-carry adders as parallel prefix adders have no advantage for these small bit widths.

The \mathbb{Z}_3 multiplier of the Mul3 module can also be directly implemented according to Equations 3 through 6 with the HPC2-SecAND gadget. The Mul3 module is fully pipelined, with a latency of five clock cycles.

Mux3 and Mux2. Mux3 can be implemented with three pipeline stages as the HPC2-SecAND gadget has a delay of two cycles for one input and one clock cycle for the other. We instantiate 13 of these two-bit MUXes in parallel in order to feed Add13 without idling.

We have a delay of two cycles for Mux2, which has two secret data input and a secret select input. We instantiate 13 MUXes in the \mathcal{R}_q multiplier to select between the CSubQ output and the non-subtracted value. We also instantiate two multiplexers during the weight check calculation to select between the original r' and the fixed vector. Finally, we use eight multiplexers to select between the encoded r' and ρ after the ciphertext comparison.

SHA-Ch and SHA-Ma. Both the SHA-Ch and SHA-Ma can be directly implemented according to Equation 11 and 12 respectively with the HPC2-SecAND gadget. We implement both operations with a width of 64 bit, in order to be able to directly feed the output to the Add64 module. The SHA-Ch has a latency of two clock cycles, while SHA-Ma has a latency of three clock cycles.

5 Evaluation

After introducing our implementation concept, we present the corresponding implementation results in this section. Furthermore, we formally verify and perform practical measurements of our building blocks in order to demonstrate their protection against side-channel attacks. Eventually, we compare our hardware implementation of Streamlined NTRU Prime to a hardware design of Saber.

5.1 Implementation Results

We implement our design on a Xilinx Artix-7 device, using Vivado v2021.2 (64-bit), for the sntrup761 parameter set. We also synthesize our design for an ASIC using the 45 nm Nangate open cell library. Table 1 shows the latency, frequency, and peak randomness demand per module and masking degree. As shown, the cycle count is dominated by the three polynomial multiplications, which take 93 % of all total cycles. At the same time, the peak randomness is always set by the 64-bit adder in the SHA-512 module. While the total cycle count is independent of the masking order, the maximum clock frequency varies: On an FPGA and at masking orders 1 and 3, the design reaches 200 MHz, but the maximum frequency is lower for masking orders 2, 4, and 5. For all three, the critical path lies in the SHA-512 module. For the ASIC, the design reaches a higher maximum clock frequency than the FPGA at first order, with 207 MHz. However, as the masking order increases, the maximum frequency drops off faster, reaching just 75 MHz at fifth order and 100 MHz at sixth order. Here, the critical path also lies in the SHA-512 module.

In Table 2, the footprint per module and masking degree is shown for Artix-7 FPGA. As expected, the area increases vastly with increasing masking degrees. Interestingly, for all masking orders, the SHA-512 dominates the resource cost consuming roughly 61 % of all LUT and FF. The next most expensive operation is the rounding during the re-encryption, followed by the \mathcal{R}_q polynomial multiplication. When comparing the ratios of cycle counts and the resources consumed, it is apparent that the current SHA-512 implementation is sub-optimal: it is too expensive when considering the whole design. In particular, the 64 bit adder is oversized. For a better ratio of cycles and resources consumed, using a smaller, e.g., 16 bit adder multiple times for each 64 bit addition, would be more efficient while adding only a comparatively minor number of cycles. Doing so would also allow the SHA-Ch and SHA-Ma gadgets to have smaller widths, saving further resources. Finally, this would reduce the maximum of random bits used per cycle.

		Maximum Randomness (bits per cycle)									
$\underline{\mathbf{Module}}$	Cycle Count				asking Oro						
		1	2	3	4	5	6	7			
Decap	1870049	52	82	156	252	370	510	672			
Encode \mathcal{R}_3	765	4	12	24	40	60	84	112			
C' comp.	4050	14	42	84	140	210	294	392			
Decrypt	1171270	96	288	576	960	1440	2016	2688			
$\mod 3$	29	46	138	276	460	690	966	1288			
Mul. \mathcal{R}_3	581 409	6	18	36	60	90	126	168			
Weight calc.	9145	42	126	252	420	630	882	1176			
Re-Encrypt	581 501	123	369	738	1230	1845	2583	3444			
Rounding	812	123	369	738	1230	1845	2583	3444			
Mul. \mathcal{R}_q	580 646	103	309	618	1030	1545	2163	2884			
Adder 13-bit	10	32	96	192	320	480	672	896			
13 Mux2	2	13	39	78	130	195	273	364			
13 Mux3	3	26	78	156	260	390	546	728			
SHA-512	7845	310	930	1860	3100	4650	6510	8 680			
SHA-Ma	2	128	384	768	1280	1920	2688	3584			
SHA-Ch	2	64	192	384	640	960	1344	1792			
Adder 64-bit	14	310	930	1860	3100	4650	6510	8680			
Total	1 870 049	310	930	1 860	3 100	4 650	6 5 1 0	8 680			
EDGA	f _{max} (MHz)	200	182	200	169	179	_	_			
FPGA	Latency (ms)	9.35	10.3	9.35	11.1	11.4	_	_			
ASIC	f _{max} (MHz)	207	165	148	91	75	100	_			
ASIC	Latency (ms)	9.03	11.3	12.6	20.5	24.9	18.7	-			

Table 1. Latency, frequency, and randomness results after Place and Route (PnR). Note that the cycle count for SHA-512 is for a single 1024-bit block. We did not perform PnR for orders 6 and 7 for an FPGA, as they no longer fit into an Artix-7 FPGA.

In the right part of Table 3, we list the gate equivalent area demand per module and masking degree for an ASIC. As we did not have access to a memory macro, we listed the memory footprint separately. We see similar behavior to the FPGA resource requirements, with the SHA-512 dominating the area footprint, followed by the rounding during the re-encryption. The total GE also grows significantly as the masking order increases, while the SRAM usage grows more slowly.

Different Masking Degrees for Decrypt and Re-Encrypt. In [ABH⁺22], the authors reason that re-encryption must be protected at a higher level than decryption during decapsulation. Our design and all building blocks can be easily adapted to any masking order allowing a flexible configuration. However, doing so would decrease the modules that can be reused across the design, e.g., the \mathcal{R}_q multiplier, which is used both during decryption and re-encryption.

5.2 Side-Channel Evaluation

In order to evaluate the protection against side-channel attacks, we rely on formal verification of each of our submodules and additionally perform practical side-channel measurements based on Test Vector Leakage Assessment (TVLA). Evaluating the entire decapsulation by formal verification or practical measurements is infeasible for typical side-channel setups due to the huge amount of required clock cycles.

Formal Verification. We formally verify the security of each module by using the recently presented verification tool VERICA [RFSG22]. VERICA is constructed based on the verification concepts developed in the side-channel analysis tool SILVER [KSM20] and the fault-injection analysis tool FIVER [RBSS⁺21]. The formal verification of a target design is performed based on its (Verilog) gate-level netlist, which is transformed into a Direct Acyclic Graph (DAG) serving as circuit model. Each node in the DAG is associated with a Binary Decision Diagram (BDD) representing the Boolean function of the corresponding

Table 2. FPGA area results after PnR. Note that this does not include the area needed for randomness generation. Not listed is the Digital Signal Processor (DSP) usage: 4 DSPs are needed as multipliers in the decoder, regardless of the masking order.

	Masking Order d											
$\underline{\mathbf{Module}}$		1			2		3			4		
	LUT	FF	BR	LUT	FF	BR	LUT	FF	BR	LUT	FF	BR
Decap	2270	1180	4.5	2493	1575	6	3088	2256	6	3766	2980	8
Encode \mathcal{R}_3	61	52	0	77	80	0	104	115	0	131	157	0
C' comp.	278	263	0	503	530	0	855	895	0	1273	1358	0
Decrypt	1743	1602	0	2680	3225	1.5	4847	5451	1.5	7276	8282	1.5
mod 3	542	719	0	1197	1528	0	2274	2638	0	3474	4049	0
Mul. \mathcal{R}_3	470	208	0	329	319	1	489	476	1	665	675	1
Weight calc.	528	612	0	1066	1286	0	1947	2194	0	2941	3350	0
Re-Encrypt	2017	2450	0.5	4138	5180	1	7755	8936	1	11696	13714	1
Rounding	1888	2387	0	4080	5108	0	7695	8851	0	11636	13616	0
Mul. \mathcal{R}_q	1846	2148	1.5	3693	4419	2	6686	7554	2	9885	11553	2.5
Adder 13-bit	627	715	0	1352	1545	0	2523	2690	0	3729	4150	0
13 Mux2	182	221	0	390	468	0	676	806	0	1040	1235	0
13 Mux3	211	286	0	463	625	0	848	1107	0	1314	1723	0
SHA-512	11684	12035	2	22493	23880	3	38370	39406	8	56207	59097	9
SHA-Ma	1528	1664	0	3439	3840	0	6624	6912	0	10197	10880	0
SHA-Ch	896	1088	0	1920	2304	0	3584	3968	0	5440	6080	0
Adder 64-bit	5996	5663	0	12352	23162	0	22506	21770	0	32702	33740	0
Total	19923	19725	8.5	36340	39209	13.5	62498	65463	18.5	91731	98726	22
Total w/o SHA	8239	7690	6.5	13847	15329	10.5	24128	26057	10.5	35524	39629	13
SHA Pct.	58.4	61.0	23.5	61.9	60.9	22.2	61.4	60.2	43.2	61.3	59.9	40.9

gate. This data structure allows efficient applications of statistical checks verifying side-channel security in the glitch-extended probing model and composability notions. To this end, we analyze our modules in the glitch-extended d-probing model for different security orders. The security order d was configured accordingly to the security order of the design under test. For the evaluation, we use a machine equipped with an Intel Xeon CPU (E5-1660) running at 3.20GHz and 128 GB of RAM. VERICA is allowed to use up to 16 cores and 8 GB of RAM per core. The corresponding results are shown in Table 4. Note that all modules pass first- and second-order verification, while third-order verification is too complex for Mod3 and Mul3. For the Add13 and Add64 modules, we use the implementation by Bache and Güneysu [BG22], which has been analyzed by practical measurements.

Measurement Results. As additional security analysis, we performed side-channel measurements of our first-order protected designs on a Sakura-G FPGA evaluation board, which is equipped with a Xilinx Spartan 6 FPGA. The target FPGA was supplied with a $4\,\mathrm{MHz}$ clock while the power consumption was measured via the voltage drop over a $1\,\Omega$ shunt resistor. The power traces were acquired by using a ZFL-2000GH+ Low Noise Amplifier (LNA) connected to a Spectrum M4 oscilloscope (8 bit resolution). The oscilloscope collected the data with a sample rate of $1.5\,\mathrm{GS/s}$. To generate the required online randomness, we instantiated a Keccak core used as Pseudorandom Number Generator (PRNG).

The measurement results for 10 million power traces can be found in Figure 4, Figure 5, and Figure 6. For all experiments, we first plot a sample trace to document a proper setup of the measurement equipment. In the subsequent plots, we used Welsh's t-test to detect potential leakage. In general and for a low number of sample points, a threshold of ± 4.5 is used to decide whether the design leaks information via the power consumption [SM15]. To this end, we do not observe any notable leakage in the first order but – as expected – some leakage in the second order.

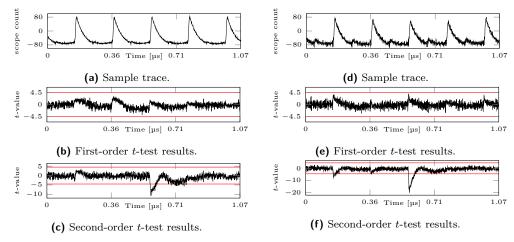


Figure 4. Measurement results for the SHA-Ma module (left) and the SHA-Ch module (right) using 10 million traces. Both modules are instantiated for d=1.

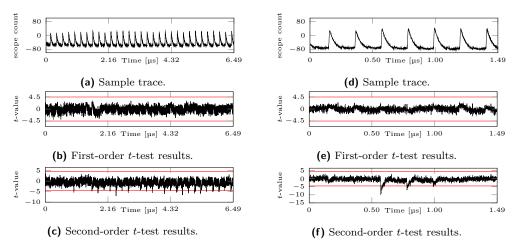


Figure 5. Measurement results for the Mod3 module (left) and the Mul3 module (right) using 10 million traces. Both modules are instantiated for d=1.

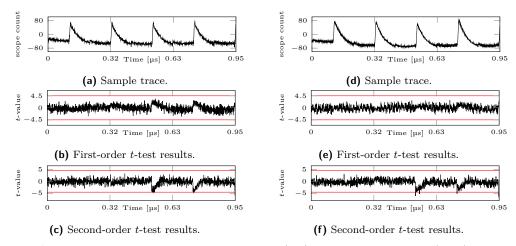


Figure 6. Measurement results for the Mux2 module (left) and the Mux3 module (right) using 10 million traces. Both modules are instantiated for d=1.

Table 3. ASIC area results in Gate Equivalent (GE), using the 45 nm Nangate open cell library. The area does not include SRAM cells, which are listed separately. Note that this does not include the area needed for randomness generation. The area for the Encode \mathcal{R}_3 entity is not available for masking orders one through three, as it was merged with its parent entity.

Module			Maskin	\mathbf{g} Order d		
Module	1	2	3	4	5	6
Decap	14 703	18 520	23 632	29 561	37 176	46 078
Encode \mathcal{R}_3	n/a	n/a	n/a	1130	1447	1799
C' comp.	2052	4103	6943	10571	14981	20208
Decrypt	14727	28021	46047	68449	95889	128306
mod 3	5744	12216	21101	32386	46085	62295
Mul. \mathcal{R}_3	2452	3436	4756	6065	8025	10412
Weight calc.	5688	11202	18584	27825	38949	51986
Re-Encrypt	29615	56009	90348	127595	176907	234225
Rounding	29057	55375	89636	126818	176050	233295
Mul. \mathcal{R}_q	25244	45906	73131	103820	143784	190601
Adder 13-bit	7115	14482	24375	36840	51817	69367
13 Mux2	1607	3510	6144	9511	13611	18442
13 Mux3	2015	4468	7910	12356	17811	24366
SHA-512	114570	218453	354242	519545	719019	950021
SHA-Ma	12416	29440	53674	85120	123776	169642
SHA-Ch	7914	17280	30250	46826	67008	90794
Adder 64-bit	55503	114820	$195\ 205$	296601	419131	563160
Total	201 112	373 349	600 100	870 124	1 204 839	1 594 022
Total w/o SHA	86542	154896	245858	350579	485820	644001
SHA Pct.	57.0	58.5	59.0	59.7	59.7	59.6
SRAM (bits)	189440	246272	294912	343296	393216	443392

Table 4. Verification results of the protected submodules using VERICA. We report for each design the number of combinational gates, memory gates and the verification time. The verification of the expected security order is indicated by green check marks. All verification results marked with ∞ are not finished in a reasonable time by VERICA.

		First Order				Second Order				Third Order			
\mathbf{Design}	comb.	mem.	sec.	time	comb.	mem.	sec.	time	comb.	mem.	sec.	time	
Mux2	16	17	1′	$0.383{\rm s}$	39	36	2	$0.402{\rm s}$	72	62	3′	20.609 s	
Mux3	28	31	1	$0.385\mathrm{s}$	72	69	2	$0.521\mathrm{s}$	136	122	3 ′	$4.985 \mathrm{h}$	
Mod3	586	774	1	$1.125\mathrm{s}$	1464	1581	2^{\checkmark}	$90.82 \mathrm{min}$	2742	2668	_	∞	
Mul3	89	94	1	$0.412\mathrm{s}$	221	204	2^{\checkmark}	$23.591\mathrm{s}$	413	356	_	∞	
SHA-Ch	16	17	1	$0.404\mathrm{s}$	39	36	2^{\checkmark}	$0.420\mathrm{s}$	72	62	3 ′	$26.355\mathrm{s}$	
SHA-Ma	28	26	1	$0.386\mathrm{s}$	72	60	2	$0.928\mathrm{s}$	136	108	3 ′	$11.5 \mathrm{h}$	

5.3 Comparison

In Table 5, we compare our implementation against an unmasked implementation of Streamlined NTRU Prime and two first-order masked FPGA implementations of Saber and Kyber. To the best of our knowledge, we are the first to report a higher-order full FPGA implementation of any PQC scheme and the first to report a masked full ASIC PQC implementation. Thus, we cannot compare it to other higher-order implementations. As expected, the two unprotected implementations are smaller, faster or both. The masked Saber implementation also has a comparable LUT and FF footprint to our first-order implementation and uses no BRAM but significantly more DSPs. However, it is about an order of magnitude faster. In contrast, the masked Kyber-512 implementation is bigger even than our fourth-order implementation, but only faster by a factor of 6.8 compared to our first-order implementation.

Moreover, both the Saber and the Kyber-512 implementations only support first order, while our design can easily be instantiated at an arbitrary level, allowing protection against more advanced attacks. Finally, our masked gadgets have been formally verified to be secure, and we do not need any masking conversion which may be used in future attacks.

Table 6 shows how performance progresses for our implementation and provides a comparison to gadget-based implementations of AES. While the randomness overhead is

	Δrea	r		may rand	
we are the first to repo	ort a fully masked	implementation of	any POC	C scheme.	
synthesized for Artix-7	7, except for Kybe	r, which is synthesiz	ed for Vi	irtex-7. Note t	that for ASIC,
Table 5. Comparison	n with other mas	sked PQC impleme	$_{ m ntations}$. All implem	entations are

Scheme	LUT	Ar FF	ea BRAM	DSP	Cycle cnt.	$\frac{f_{max}}{\text{MHz}}$	max rand. bits / cycle	d	Ref.
sNTRUp-761 sNTRUp-761	36 789 6 279	22 700 3 086	3.5 3.0	9 7	$10989 \\ 85628$	137 131	0	0	[PMT ⁺ 22] [PMT ⁺ 22]
sNTRUp-761 sNTRUp-761 sNTRUp-761 sNTRUp-761 Saber Kyber-512	19 923 36 340 62 498 91 731 19 299 152 860	19 725 39 209 65 463 98 726 21 977 DNR	8.5 13.5 18.5 22.0 0.0 489.5	4 4 4 4 64 76	1870 049 1870 049 1870 049 1870 049 72 005 137 738	200 182 200 169 125 100	310 930 1 860 3 100 DNR DNR	1 2 3 4 1	this this this this [AMD ⁺ 21] [KNAH22]

Table 6. Comparison with gadget-based masked implementations of symmetric schemes. Overhead is given as the fraction between the current row and the previous row minus one.

			Utilization				ing	
		Area	overhead	Rand.	overhead	Latency	overhead	
Scheme	d	[GE]	[%]	[bit]	[%]	[ns]	[%]	Ref.
	1	201 112	_	310	_	9.05×10^{6}	_	
sNTRUp-761	2	373349	85.6	930	200	11.3×10^{6}	25.9	${f this}$
	3	600100	60.7	1860	100	12.6×10^{6}	11.5	
	0	3 263	_	0	_	189.2	_	
AES (serial)	1	10090	209	34	∞	4311	2188	[KMMS22]
ALS (serial)	2	17649	74.9	102	200	5434	26.1	[KWW522]
	3	27026	53.1	204	100	5537	1.88	
	0	9 906	_	0	_	20.35	_	
AES (round-b.)	1	52597	431	680	∞	201.9	892	[IZMMC99]
	2	131631	150	2040	200	236.6	17.2	[KMMS22]
	3	246924	87.6	4080	100	250.5	5.86	

independent of the scheme (it only depends on the gadget, which is HPC2 for all schemes in this table), differences can be observed for area and delay. Still, the relative area overhead is similar for Streamlined NTRU Prime and AES. However, the delay overhead is slightly worse for d=3. The reason for this is likely the larger absolute area of Streamlined NTRU Prime, which causes additional routing delays.

6 Discussion

This section addresses and discusses potential improvements and the huge overhead introduced by masking the symmetric core in Streamlined NTRU Prime. Additionally, we briefly discuss applying our concepts and approaches to Kyber.

6.1 Gadget-based Masking

There are several advantages in a gadget-based masked implementation. First, it is effortless to adapt to an arbitrary masking order. This obviously reduces the time required for the development. Moreover, no masking conversion can be attacked since there is none. The masking conversion was the target in the attacks against a first-order and third-order masked Saber implementation [NDJ21,NWDP22]. Additionally, exchanging the underlying gadgets with others with the same latency properties is usually straightforward. For example, it could be possible to achieve a fault-secure implementation easily by deploying the work from [FRBSG22]. In addition, while our design does not include an RNG,

generating randomness is relatively straightforward and cheap in hardware: The recent work [CMM $^+23$] analyses the cost of securely generating random bits for use in masking, with costs of 20 to 30 GE or 3 to 4 LUTs per bit while using a round-reduced version of the Trivium cipher. This additional area is minimal compared to the area usage of our design, and would only add roughly 4%.

6.2 Potential Improvements

We leave several potential improvements as future work and address them here. The polynomial multiplications have the most conspicuous latencies, where the two \mathcal{R}_q multiplications take 62% of the decapsulation cycle counts, and the multiplication in \mathcal{R}_3 takes another 31%. To speed this up, it is possible to instantiate more adders in parallel at the cost of slightly more area and a potentially higher amount of randomness per clock cycle, depending on the grade of deployed parallelism. Thus, halving the latency of both multipliers results in a 47% speed-up at the cost of approximately 8% more gate equivalents for the first-order ASIC implementation. Moreover, a potential area reduction can be achieved by optimizing the CSubQ module including a positive impact on the amount of required randomness.

Additionally, we want to stress that the specified encoding procedure for polynomials \mathcal{R}_q is suboptimal for hardware implementations, as it includes multiplications. This accounts for the four DSP slices required in the FPGA implementation and about 7.3 kGE in the ASIC implementation. However, alternatives would increase transmission sizes and obviously need a change of the Streamlined NTRU Prime specification.

6.3 Symmetric Core

As discussed in Section 5.1, masking the symmetric core (i.e., SHA-512) in Streamlined NTRU Prime consumes a considerable large part of the entire implementation's footprint and has the highest per cycle randomness consumption. It also limits the maximum frequency due to the high routing cost of the 64-bit Sklansky adder. Nevertheless, hardened SHA-512 implementations are widely deployed in industry and can, for example, be found in smartcards and secure elements. Thus, one could assume that a secure SHA-512 is already available and does not need to be implemented. If we exclude the SHA-512 from the area consumption (cf. Table 2 and Table 3), then the design is not only surprisingly small at first order, but the area overhead is much more moderate with increasing masking order.

Another possibility would be to replace the 64-bit Sklansky adder deployed in the SHA-512 module by a smaller one, trading area for latency. Moreover, it is possible to deploy no additional adder for the SHA-512 module by reusing the secure adder from the polynomial multiplication module. In this case, five consecutive 13-bit additions would yield the 64-bit addition. This would require cleverly scheduling the additions required by SHA-512 such that the 13-bit adder pipeline is maximally occupied. As can be seen from Table 2 and Table 3, the 64-bit Sklansky adder occupies about half of the area of the SHA-512 module and about a quarter of the overall area.

Additionally, in order to reduce the total area overhead introduced by the masked symmetric core in Streamlined NTRU Prime, the SHA-512 could be replaced by an implementation based on Keccak [BDPA13]. As Keccak does not use an adder internally, it is significantly easier and cheaper to mask. Most notably, it can be implemented with a very low amount of fresh randomness [BDN $^+$ 13]. In addition, as the critical path lies in the SHA-512 module for both FPGAs and ASICs, using Keccak would likely increase the maximum achievable clock frequency. However, this would deviate from the Streamlined NTRU Prime specification and would not be interoperable with other Streamlined NTRU Prime implementations.

			1	, , , , , , , , , , , , , , , , , , ,	
Scheme	NIST Category	Polynomial size	Module size		of additions small coefficients
Kyber-512 sNTRUp	I	256 653	k = 2	$527104 \\ 854124$	$0 \\ 427062$
sNTRUp	II	761	_	1 159 764	579 882
Kyber-768 sNTRUp	III	256 857	k = 3	988160 1470612	0 735 306
sNTRUp sNTRUp	IV IV	953 1013	_	$1818324 \\ 2054364$	909 162 1 027 182
Kyber-1024 sNTRUp	V	256 1277	k = 4	1 580 800 3 264 012	0 1 632 006

Table 7. Comparison to Kyber

6.4 Applicability to Kyber

The efficiency of our gadget-based masking is built upon the fact that the three polynomial multiplications that are carried out each include a secret polynomial with ternary coefficients, where the other one is either small and secret as well or has a big coefficient modulus and is public. This enables us to perform schoolbook multiplication in Boolean domain. Notably, Kyber has a similar property: Here, all polynomial multiplications have one public input polynomial with "big" coefficients modulo $q=3\,329$.

Moreover, the polynomial degree is far smaller, with 256 compared to 761 for Streamlined NTRU Prime, enabling a faster multiplication. For Kyber, $256^2=65\,536$ coefficient additions are to be performed per polynomial multiplication, whereas Streamlined NTRU Prime with p=761 requires $p^2+p=579\,882$ coefficient additions. However, Kyber requires more multiplications to be performed: for $k\in\{2,3,4\}$, it requires k^2+2k polynomial multiplications, as well as k^2+4k-1 polynomial additions, whereas Streamlined NTRU Prime constantly requires three polynomial multiplications, one of which only uses "small" coefficients and is thus much cheaper.

We compare the cost in terms of the estimated number of coefficient additions in Table 7. As seen there, Kyber consistently requires fewer "big" coefficient additions than Streamlined NTRU Prime in the regarding security categories. Another advantage for Kyber is that during key generation, it features no operations that are infeasible to mask in Boolean domain, which is in contrast to Streamlined NTRU Prime, where this is not possible. The most complex remaining operations in Kyber, both for key generation and decapsulation, are (de-)compression and sampling for a centered binomial distribution using a Keccak output stream, both of which are feasible in Boolean domain.

One downside for Kyber is that the secret coefficients have the range of [-2,2] or [-3,3]. This would require a more complex five-way or seven-way secure multiplexer. In addition, a gadget-based masked Kyber implementation would require a Number-Theoretic Transform (NTT) core: Kyber requires extending a seed into a public matrix of polynomials, which are assumed to be in NTT domain. Since the implementation would not perform multiplication in NTT domain, an inverse transform of each polynomial in the matrix would be required, resulting in k^2 inverse NTTs during decapsulation. Finally, it is noteworthy that the fact that Kyber uses the same polynomial ring for all security levels is no advantage for a gadget-based masked implementation since schoolbook multiplication is used for Streamlined NTRU Prime, which also allows for easy parametrization. On the other hand, Streamlined NTRU Prime changes the coefficient modulus over the parameter sets, which might require manual adjustments. Overall, we leave this as an interesting open idea for future work.

6.5 Applicability to Dilithium

From a side-channel point of view, multiplying the public challenge polynomial c with the secret key vector $\mathbf{s_1}$, followed by an addition to the nonce \mathbf{y} , is the most critical operation for signature generation in Dilithium [SLKG23]. We highlight that the challenge polynomial c is also public for rejected signature candidates [KLRBG23] and is sparse and ternary. The coefficients of the secret key polynomial vector $\mathbf{s_1}$ are uniformly random with a small bound (i.e., there are only five or nine possible values). The nonce \mathbf{y} , finally, is a secret with a large bound. Thus, the whole operation can be masked similarly to the technique explained in this paper but with sparse multiplication. Notably, no modular reduction for the coefficient accumulation is required as the maximum bound for each coefficient after the operation is lower than the modulus, and the subsequent bound check can be performed in signed representation. The multiplication-accumulation $\mathbf{w} - c\mathbf{s_2}$ of the secret \mathbf{w} and $\mathbf{s_2}$ and public c can also be masked using our method, though it does require modular reduction.

On the other hand, the other critical multiplication during signature generation - $\mathbf{A}\mathbf{y}$ – fulfills only the criterion of having one public (\mathbf{A}) and one secret factor (\mathbf{y}). The coefficients of \mathbf{y} have 2^{18} or 2^{20} possible values (depending on the parameter set), rendering the gadget-based masking approach infeasible for this operation as a full multiplier is required. In addition, similar to Kyber, \mathbf{A} is extended from seeds and expected to be in the NTT domain.

7 Conclusion

In our work, we have presented the first gadget-based masked implementation of any PKC scheme. Notably, it is competitive regarding area demand to other protected PQC implementations while offering reasonable latency. The main advantage is the ability to adapt the implementation easily to arbitrary masking orders. For the first-order secure instance of the implementation, 19 923 LUTs, 19 725 FFs, and 8.5 BRAMs are utilized, reaching a frequency of 200 MHz. Implemented as an ASIC, the first-order secure instance consumes 201k GE and 189 kbit SRAM, reaching a frequency of 207 MHz. This results in a latency of only 9.35 ms on an FPGA and 9.03 ms on an ASIC, with a peak demand of fresh randomness of 310 bit per clock cycle. While for higher masking degrees, the latency only increases slightly due to a lower frequency, the randomness demand increases to 3 100 bit per clock cycle for d=4. Nevertheless, further optimization of the hashing module could significantly reduce the area and randomness consumption. Moreover, our first-order implementation is formally and practically verified to be secure. Finally, we also analyzed the applicability of our concept to the designated NIST standard algorithm Kyber, finding that gadget-based masking could be applied efficiently as well.

Acknowledgments

We would like to thank Joppe Bos, Daniel J. Bernstein, Bo-Yuan Peng and Bo-Yin Yang for their help. This work was labelled by the EUREKA cluster PENTA and funded by German authorities under grant agreement PENTA-2018e-17004-SunRISE. This work was supported by the Federal Ministry of Education and Research (BMBF) of the Federal Republic of Germany (grants 16KIS1572K, SASVI and 16KISK038, 6GEM). The work described in this paper has been supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy - EXC 2092 CASA - 390781972, and by the European Commission under the grant agreement number 101070374 (CONVOLVE).

References

- [ABH+22] Melissa Azouaoui, Olivier Bronchain, Clément Hoffmann, Yulia Kuzovkova, Tobias Schneider, and François-Xavier Standaert. Systematic Study of Decryption and Re-encryption Leakage: The Case of Kyber. In Josep Balasch and Colin O'Flynn, editors, Constructive Side-Channel Analysis and Secure Design 13th International Workshop, COSADE 2022, Leuven, Belgium, April 11-12, 2022, Proceedings, volume 13211 of Lecture Notes in Computer Science, pages 236–256. Springer, 2022.
- [ACC⁺21] Erdem Alkim, Dean Yun-Li Cheng, Chi-Ming Marvin Chung, Hülya Evkan, Leo Wei-Lun Huang, Vincent Hwang, Ching-Lin Trista Li, Ruben Niederhagen, Cheng-Jhih Shih, Julian Wälde, and Bo-Yin Yang. Polynomial multiplication in NTRU prime. *IACR TCHES*, 2021(1):217–238, 2021. https://tches.iacr.org/index.php/TCHES/article/view/8733.
- [AMD+21] Abubakr Abdulgadir, Kamyar Mohajerani, Viet Ba Dang, Jens-Peter Kaps, and Kris Gaj. A Lightweight Implementation of Saber Resistant Against Side-Channel Attacks. In Avishek Adhikari, Ralf Küsters, and Bart Preneel, editors, Progress in Cryptology INDOCRYPT 2021 22nd International Conference on Cryptology in India, Jaipur, India, December 12-15, 2021, Proceedings, volume 13143 of Lecture Notes in Computer Science, pages 224–245. Springer, 2021.
- [AR21] Amund Askeland and Sondre Rønjom. A side-channel assisted attack on NTRU. Cryptology ePrint Archive, Report 2021/790, 2021. https://eprint.iacr.org/2021/790.
- [BBC⁺20] Daniel J. Bernstein, Billy Bob Brumley, Ming-Shing Chen, Chitchanok Chuengsatiansup, Tanja Lange, Adrian Marotzke, Bo-Yuan Peng, Nicola Tuveri, Christine van Vredendaal, and Bo-Yin Yang. NTRU Prime. Technical report, National Institute of Standards and Technology, 2020. available at https://csrc.nist.gov/projects/post-quantum-cryptography-standardization/round-3-submissions.
- [BBD+15] Gilles Barthe, Sonia Belaïd, François Dupressoir, Pierre-Alain Fouque, Benjamin Grégoire, and Pierre-Yves Strub. Verified proofs of higher-order masking. In Elisabeth Oswald and Marc Fischlin, editors, EUROCRYPT 2015, Part I, volume 9056 of LNCS, pages 457–485. Springer, Heidelberg, April 2015.
- [BBD⁺16] Gilles Barthe, Sonia Belaïd, François Dupressoir, Pierre-Alain Fouque, Benjamin Grégoire, Pierre-Yves Strub, and Rébecca Zucchini. Strong non-interference and type-directed higher-order masking. In Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi, editors, ACM CCS 2016, pages 116–129. ACM Press, October 2016.
- [BCLv17] Daniel J. Bernstein, Chitchanok Chuengsatiansup, Tanja Lange, and Christine van Vredendaal. NTRU prime: Reducing attack surface at low cost. In Carlisle Adams and Jan Camenisch, editors, SAC 2017, volume 10719 of LNCS, pages 235–260. Springer, Heidelberg, August 2017.
- [BDN⁺13] Begül Bilgin, Joan Daemen, Ventzislav Nikov, Svetla Nikova, Vincent Rijmen, and Gilles Van Assche. Efficient and First-Order DPA Resistant Implementations of Keccak. In Aurélien Francillon and Pankaj Rohatgi, editors, Smart Card Research and Advanced Applications 12th International Conference, CARDIS 2013, Berlin, Germany, November 27-29, 2013. Revised Selected

- Papers, volume 8419 of Lecture Notes in Computer Science, pages 187–199. Springer, 2013.
- [BDPA13] Guido Bertoni, Joan Daemen, Michaël Peeters, and Gilles Van Assche. Keccak. In Thomas Johansson and Phong Q. Nguyen, editors, Advances in Cryptology - EUROCRYPT 2013, 32nd Annual International Conference on the Theory and Applications of Cryptographic Techniques, Athens, Greece, May 26-30, 2013. Proceedings, volume 7881 of Lecture Notes in Computer Science, pages 313-314. Springer, 2013.
- [BG22] Florian Bache and Tim Güneysu. Boolean Masking for Arithmetic Additions at Arbitrary Order in Hardware. *Applied Sciences*, 12(5):2274, 2022.
- [BGI⁺18] Roderick Bloem, Hannes Groß, Rinat Iusupov, Bettina Könighofer, Stefan Mangard, and Johannes Winter. Formal verification of masked hardware implementations in the presence of glitches. In Jesper Buus Nielsen and Vincent Rijmen, editors, EUROCRYPT 2018, Part II, volume 10821 of LNCS, pages 321–353. Springer, Heidelberg, April / May 2018.
- [CGLS21] Gaëtan Cassiers, Benjamin Grégoire, Itamar Levi, and François-Xavier Standaert. Hardware Private Circuits: From Trivial Composition to Full Verification. IEEE Trans. Computers, 70(10):1677–1690, 2021.
- [CGTZ23] Jean-Sébastien Coron, François Gérard, Matthias Trannoy, and Rina Zeitoun. High-order masking of NTRU. IACR Trans. Cryptogr. Hardw. Embed. Syst., 2023(2):180–211, 2023.
- [CHK+21] Chi-Ming Marvin Chung, Vincent Hwang, Matthias J. Kannwischer, Gregor Seiler, Cheng-Jhih Shih, and Bo-Yin Yang. NTT multiplication for NTT-unfriendly rings. *IACR TCHES*, 2021(2):159–188, 2021. https://tches.iacr.org/index.php/TCHES/article/view/8791.
- [CMM⁺23] Gaëtan Cassiers, Loïc Masure, Charles Momin, Thorben Moos, Amir Moradi, and François-Xavier Standaert. Randomness generation for secure hardware masking-unrolled trivium to the rescue. *Cryptology ePrint Archive*, 2023.
- [CS20] Gaëtan Cassiers and François-Xavier Standaert. Trivially and Efficiently Composing Masked Gadgets With Probe Isolating Non-Interference. *IEEE Trans. Inf. Forensics Secur.*, 15:2542–2555, 2020.
- [DGBN09] Jean-Luc Danger, Sylvain Guilley, Shivam Bhasin, and Maxime Nassar. Overview of dual rail with precharge logic styles to thwart implementation-level attacks on hardware cryptoprocessors. In 2009 3rd International Conference on Signals, Circuits and Systems (SCS), pages 1–8. IEEE, 2009.
- [DNG22] Elena Dubrova, Kalle Ngo, and Joel Gärtner. Breaking a fifth-order masked implementation of crystals-kyber by copy-paste. *Cryptology ePrint Archive*, 2022.
- [FBR+22] Tim Fritzmann, Michiel Van Beirendonck, Debapriya Basu Roy, Patrick Karl, Thomas Schamberger, Ingrid Verbauwhede, and Georg Sigl. Masked Accelerators and Instruction Set Extensions for Post-Quantum Cryptography. IACR Trans. Cryptogr. Hardw. Embed. Syst., 2022(1):414-460, 2022.
- [FGP⁺18] Sebastian Faust, Vincent Grosso, Santos Merino Del Pozo, Clara Paglialonga, and François-Xavier Standaert. Composable masking schemes in the presence

- of physical defaults & the robust probing model. *IACR TCHES*, 2018(3):89–120, 2018. https://tches.iacr.org/index.php/TCHES/article/view/7270.
- [FRBSG22] Jakob Feldtkeller, Jan Richter-Brockmann, Pascal Sasdrich, and Tim Güneysu. CINI MINIS: Domain isolation for fault and combined security. In Heng Yin, Angelos Stavrou, Cas Cremers, and Elaine Shi, editors, *ACM CCS 2022*, pages 1023–1036. ACM Press, November 2022.
- [HHP+21] Mike Hamburg, Julius Hermelink, Robert Primas, Simona Samardjiska, Thomas Schamberger, Silvan Streit, Emanuele Strieder, and Christine van Vredendaal. Chosen ciphertext k-trace attacks on masked CCA2 secure kyber.

 IACR TCHES, 2021(4):88-113, 2021. https://tches.iacr.org/index.php/TCHES/article/view/9061.
- [HPP21] Julius Hermelink, Peter Pessl, and Thomas Pöppelmann. Fault-Enabled Chosen-Ciphertext Attacks on Kyber. In Avishek Adhikari, Ralf Küsters, and Bart Preneel, editors, Progress in Cryptology INDOCRYPT 2021 22nd International Conference on Cryptology in India, Jaipur, India, December 12-15, 2021, Proceedings, volume 13143 of Lecture Notes in Computer Science, pages 311–334. Springer, 2021.
- [ISW03] Yuval Ishai, Amit Sahai, and David Wagner. Private circuits: Securing hardware against probing attacks. In Dan Boneh, editor, *CRYPTO 2003*, volume 2729 of *LNCS*, pages 463–481. Springer, Heidelberg, August 2003.
- [KA21] Emre Karabulut and Aydin Aysu. FALCON Down: Breaking FALCON Post-Quantum Signature Scheme through Side-Channel Attacks. In 58th ACM/IEEE Design Automation Conference, DAC 2021, San Francisco, CA, USA, December 5-9, 2021, pages 691–696. IEEE, 2021.
- [KAA21] Emre Karabulut, Erdem Alkim, and Aydin Aysu. Single-Trace Side-Channel Attacks on ω -Small Polynomial Sampling: With Applications to NTRU, NTRU Prime, and CRYSTALS-DILITHIUM. In *IEEE International Symposium on Hardware Oriented Security and Trust, HOST 2021, Tysons Corner, VA, USA, December 12-15, 2021*, pages 35–45. IEEE, 2021.
- [KLRBG23] Markus Krausz, Georg Land, Jan Richter-Brockmann, and Tim Güneysu. A holistic approach towards side-channel secure fixed-weight polynomial sampling. In Alexandra Boldyreva and Vladimir Kolesnikov, editors, PKC 2023, Part II, volume 13941 of LNCS, pages 94–124. Springer, Heidelberg, May 2023.
- [KM22] David Knichel and Amir Moradi. Low-latency hardware private circuits. In Heng Yin, Angelos Stavrou, Cas Cremers, and Elaine Shi, editors, ACM CCS 2022, pages 1799–1812. ACM Press, November 2022.
- [KMMS22] David Knichel, Amir Moradi, Nicolai Müller, and Pascal Sasdrich. Automated generation of masked hardware. *IACR TCHES*, 2022(1):589–629, 2022.
- [KNAH22] Tendayi Kamucheka, Alexander Nelson, David Andrews, and Miaoqing Huang. A masked pure-hardware implementation of kyber cryptographic algorithm. In International Conference on Field-Programmable Technology, (IC)FPT 2022, Hong Kong, December 5-9, 2022, page 1. IEEE, 2022.

- [KSM20] David Knichel, Pascal Sasdrich, and Amir Moradi. SILVER statistical independence and leakage verification. In Shiho Moriai and Huaxiong Wang, editors, ASIACRYPT 2020, Part I, volume 12491 of LNCS, pages 787–816. Springer, Heidelberg, December 2020.
- [KSM22] David Knichel, Pascal Sasdrich, and Amir Moradi. Generic Hardware Private Circuits Towards Automated Generation of Composable Secure Gadgets. IACR Trans. Cryptogr. Hardw. Embed. Syst., 2022(1):323–344, 2022.
- [Mar20] Adrian Marotzke. A Constant Time Full Hardware Implementation of Stream-lined NTRU Prime. In Pierre-Yvan Liardet and Nele Mentens, editors, Smart Card Research and Advanced Applications 19th International Conference, CARDIS 2020, Virtual Event, November 18-19, 2020, Revised Selected Papers, volume 12609 of Lecture Notes in Computer Science, pages 3–17. Springer, 2020.
- [MKEP11] Amir Moradi, Mario Kirschbaum, Thomas Eisenbarth, and Christof Paar. Masked dual-rail precharge logic encounters state-of-the-art power analysis methods. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 20(9):1578–1589, 2011.
- [NDGJ21] Kalle Ngo, Elena Dubrova, Qian Guo, and Thomas Johansson. A side-channel attack on a masked IND-CCA secure saber KEM implementation. *IACR TCHES*, 2021(4):676-707, 2021. https://tches.iacr.org/index.php/TCH ES/article/view/9079.
- [NDJ21] Kalle Ngo, Elena Dubrova, and Thomas Johansson. Breaking masked and shuffled CCA secure saber KEM by power analysis. In Chip-Hong Chang, Ulrich Rührmair, Stefan Katzenbeisser, and Debdeep Mukhopadhyay, editors, ASHES@CCS 2021: Proceedings of the 5th Workshop on Attacks and Solutions in Hardware Security, Virtual Event, Republic of Korea, 19 November 2021, pages 51–61. ACM, 2021.
- [NWDP22] Kalle Ngo, Ruize Wang, Elena Dubrova, and Nils Paulsrud. Side-channel attacks on lattice-based KEMs are not prevented by higher-order masking. Cryptology ePrint Archive, Report 2022/919, 2022. https://eprint.iacr.org/2022/919.
- [PKZM07] Thomas Popp, Mario Kirschbaum, Thomas Zefferer, and Stefan Mangard. Evaluation of the masked logic style MDPL on a prototype chip. In Pascal Paillier and Ingrid Verbauwhede, editors, *CHES 2007*, volume 4727 of *LNCS*, pages 81–94. Springer, Heidelberg, September 2007.
- [PMT⁺22] Bo-Yuan Peng, Adrian Marotzke, Ming-Han Tsai, Bo-Yin Yang, and Ho-Lin Chen. Streamlined NTRU Prime on FPGA. *Journal of Cryptographic Engineering*, pages 1–20, 2022.
- [RBSS⁺21] Jan Richter-Brockmann, Aein Rezaei Shahmirzadi, Pascal Sasdrich, Amir Moradi, and Tim Güneysu. FIVER robust verification of countermeasures against fault injections. *IACR TCHES*, 2021(4):447–473, 2021. https://tches.iacr.org/index.php/TCHES/article/view/9072.
- [RFSG22] Jan Richter-Brockmann, Jakob Feldtkeller, Pascal Sasdrich, and Tim Güneysu. VERICA - Verification of Combined Attacks Automated formal verification of security against simultaneous information leakage and tampering. IACR Trans. Cryptogr. Hardw. Embed. Syst., 2022(4):255–284, 2022.

- [RRCB20] Prasanna Ravi, Sujoy Sinha Roy, Anupam Chattopadhyay, and Shivam Bhasin. Generic side-channel attacks on CCA-secure lattice-based PKE and KEMs. *IACR TCHES*, 2020(3):307–335, 2020. https://tches.iacr.org/index.php/TCHES/article/view/8592.
- [SGD+09] Laurent Sauvage, Sylvain Guilley, Jean-Luc Danger, Yves Mathieu, and Maxime Nassar. Successful attack on an fpga-based wddl des cryptoprocessor without place and route constraints. In 2009 Design, Automation & Test in Europe Conference & Exhibition, pages 640–645. IEEE, 2009.
- [SLKG23] Hauke Malte Steffen, Georg Land, Lucie Johanna Kogelheide, and Tim Güneysu. Breaking and protecting the crystal: Side-channel analysis of dilithium in hardware. In Thomas Johansson and Daniel Smith-Tone, editors, Post-Quantum Cryptography 14th International Workshop, PQCrypto 2023, College Park, MD, USA, August 16-18, 2023, Proceedings, volume 14154 of Lecture Notes in Computer Science, pages 688-711. Springer, 2023.
- [SM15] Tobias Schneider and Amir Moradi. Leakage assessment methodology A clear roadmap for side-channel evaluations. In Tim Güneysu and Helena Handschuh, editors, *CHES 2015*, volume 9293 of *LNCS*, pages 495–513. Springer, Heidelberg, September 2015.
- [SMG15] Tobias Schneider, Amir Moradi, and Tim Güneysu. Arithmetic addition over Boolean masking - towards first- and second-order resistance in hardware. In Tal Malkin, Vladimir Kolesnikov, Allison Bishop Lewko, and Michalis Polychronakis, editors, ACNS 15, volume 9092 of LNCS, pages 559–578. Springer, Heidelberg, June 2015.
- [SPH22] Bo-Yeon Sim, Aesun Park, and Dong-Guk Han. Chosen-ciphertext clustering attack on CRYSTALS-KYBER using the side-channel leakage of barrett reduction. *IEEE Internet Things J.*, 9(21):21382–21397, 2022.
- [TV04] Kris Tiri and Ingrid Verbauwhede. A logic level design methodology for a secure dpa resistant asic or fpga implementation. In *Proceedings Design*, Automation and Test in Europe Conference and Exhibition, volume 1, pages 246–251. IEEE, 2004.
- [XPR⁺21] Zhuang Xu, Owen Michael Pemberton, Sujoy Sinha Roy, David Oswald, Wang Yao, and Zhiming Zheng. Magnifying side-channel leakage of lattice-based cryptosystems with chosen ciphertexts: The case study of kyber. *IEEE Transactions on Computers*, 2021.