

# RoK, Paper, SISsors – Toolkit for Lattice-based Succinct Arguments

(Full Version)

Michael Kloof<sup>1\*</sup>, Russell W. F. Lai<sup>2</sup>, Ngoc Khanh Nguyen<sup>3</sup>, and Michał Osadnik<sup>2</sup>

<sup>1</sup> ETH Zurich, Zurich, Switzerland

<sup>2</sup> Aalto University, Espoo, Finland

<sup>3</sup> King’s College London, London, UK

**Abstract.** Lattice-based succinct arguments allow to prove bounded-norm satisfiability of relations, such as  $f(\mathbf{s}) = \mathbf{t} \bmod q$  and  $\|\mathbf{s}\| \leq \beta$ , over specific cyclotomic rings  $\mathcal{O}_{\mathcal{K}}$ , with proof size polylogarithmic in the witness size. However, state-of-the-art protocols require either 1) a super-polynomial size modulus  $q$  due to a soundness gap in the security argument, or 2) a verifier which runs in time linear in the witness size. Furthermore, construction techniques often rely on specific choices of  $\mathcal{K}$  which are not mutually compatible. In this work, we exhibit a diverse toolkit for constructing efficient lattice-based succinct arguments:

- (i) We identify new subtractive sets for general cyclotomic fields  $\mathcal{K}$  and their maximal real subfields  $\mathcal{K}^+$ , which are useful as challenge sets, e.g. in arguments for exact norm bounds.
- (ii) We construct modular, verifier-succinct reductions of knowledge for the bounded-norm satisfiability of structured-linear/inner-product relations, without any soundness gap, under the vanishing SIS assumption, over any  $\mathcal{K}$  which admits polynomial-size subtractive sets.
- (iii) We propose a framework to use twisted trace maps, i.e. maps of the form  $\tau(z) = \frac{1}{N} \cdot \text{Trace}_{\mathcal{K}/\mathbb{Q}}(\alpha \cdot z)$ , to embed  $\mathbb{Z}$ -inner-products as  $\mathcal{R}$ -inner-products for some structured subrings  $\mathcal{R} \subseteq \mathcal{O}_{\mathcal{K}}$  whenever the conductor has a square-free odd part.
- (iv) We present a simple extension of our reductions of knowledge for proving the consistency between the coefficient embedding and the Chinese Remainder Transform (CRT) encoding of  $\mathbf{s}$  over any cyclotomic field  $\mathcal{K}$  with a smooth conductor, based on a succinct decomposition of the CRT map into automorphisms, and a new, simple succinct argument for proving automorphism relations.

Combining all techniques, we obtain, for example, verifier-succinct arguments for proving that  $\mathbf{s}$  satisfying  $f(\mathbf{s}) = \mathbf{t} \bmod q$  has binary coefficients, without soundness gap and with polynomial-size modulus  $q$ .

## 1 Introduction

A fundamental and recurring task in constructing lattice-based succinct arguments is to prove knowledge of a committed vector  $\mathbf{s} \in \mathcal{R}^m$  over a ring  $\mathcal{R}$  which satisfies norm-bound constraints, such as  $\|\mathbf{s}\| \leq \beta$ . For instance, such protocols could be extended directly into a succinct argument for structured languages [CLM23], combined with quadratic functional commitments to yield succinct arguments for NP [ACL<sup>+</sup>22, CLM23]<sup>4</sup>, or transformed into polynomial commitment schemes [FMN23, AFLN23, CMNW24] which allow compiling polynomial interactive oracle proofs [BCS16] into succinct arguments.

As evidenced in prior works [Lyu12, LNP22, BS23], the currently most efficient lattice-based (non-)succinct arguments operate over rings of integers  $\mathcal{R} := \mathbb{Z}[\zeta]$  of cyclotomic number fields  $\mathcal{K} := \mathbb{Q}(\zeta)$ , where  $\zeta$  is a primitive  $f$ -th root of unity for  $f = \text{poly}(\lambda)$ . Indeed, the ability to construct exponential-sized low-norm challenge sets over  $\mathcal{R}$  allows the aforementioned protocols to achieve negligible soundness in one-shot while maintaining relatively small lattice parameters. However, this comes at a cost of the following two complications.

---

\*Work done at Aalto University. The author’s affiliation changed before publication.

<sup>4</sup>[ACL<sup>+</sup>22, CLM23] relied on the knowledge-kRISIS assumption for the knowledge soundness of well-formedness of commitments. However, the assumption has subsequently been cryptanalysed [WW23, DAFS24], rendering the security proofs vacuous.

**Correctness Gap.** The first one can be described as the *correctness gap*. Namely, most of the recursion-based protocols start with the initial witness  $\mathbf{s}_0 := \mathbf{s}$ , and in the  $i$ -th iteration, an honest prover somehow folds the “current” witness  $\mathbf{s}_{i-1}$  into a new one  $\mathbf{s}_i$ ; thus shrinking the dimension of the witness, but simultaneously, increasing its norm. At the end, say after  $\mu$  iterations, the prover outputs the final witness  $\mathbf{s}_\mu$  of small (potentially constant) dimension. Suppose there exists some  $\gamma$  such that for all  $i = 1, \dots, \mu$  we have  $\|\mathbf{s}_i\| \leq \gamma \cdot \|\mathbf{s}_{i-1}\|$ . Then, in order to maintain correctness, one must inherently choose  $q > \gamma^\mu \cdot \beta \geq \|\mathbf{s}_\mu\|$ . We call this phenomenon the correctness gap, since if our only task were to commit to  $\mathbf{s}$  using a standard lattice-based commitment scheme, setting  $q = O(\beta)$  would suffice<sup>5</sup>.

**Soundness Gap.** A more concerning issue is the *soundness gap*. A vast majority of prior works based on cyclotomic rings encounter the problem that the extracted witness  $\bar{\mathbf{s}}$  is not necessarily short, but it is of the fractional form  $\bar{\mathbf{s}} := \bar{\mathbf{z}}/\bar{c} \pmod q$ , where  $q$  is the proof system modulus and both  $\bar{\mathbf{z}} \in \mathcal{R}^m$  and  $\bar{c} \in \mathcal{R}$  are somewhat short (but  $\|\bar{\mathbf{z}}\|$  is larger than  $\beta$ ). Even though this *relaxed* soundness suffices to construct basic primitives, such as signature schemes [Lyu12, DKL<sup>+</sup>18], verifiable encryption [LN17], or few-time verifiable random functions [EKS<sup>+</sup>21], it is not enough when the required functionality naturally involves proving exact norm bounds (e.g. in set membership and range proofs). But especially in the context of succinct arguments built in a recursive manner, dealing with the slack and other norm-growth related issues have shown to have enormous impact on setting up the parameters [BLNS20, BCS23, AL21, AFLN23], such as picking super-polynomial modulus  $q$ , which makes the aforementioned schemes seem barely practical.

**Prior works.** Since the soundness gap seemed to be the main efficiency bottleneck of lattice-based succinct arguments, several works naturally tried to address this issue first. To begin with, Albrecht and Lai [AL21] designed a lattice-based argument of polylogarithmic size, where the extracted witness  $\bar{\mathbf{s}}$  is somewhat short. The key ingredient of [AL21] was the notion of *subtractive sets*. Namely, a set  $S \subseteq \mathcal{R}$  is called subtractive if for any two distinct elements  $c, c' \in S$ ,  $c - c'$  is invertible over the ring  $\mathcal{R}$ . Since the invertibility is independent of the proof system modulus  $q$ , the latter can be picked freely so that the inverse  $(c - c')^{-1}$  is short relative to  $q$ . Further, it was shown how to construct such subtractive sets of cardinality  $p$  in cyclotomic rings of prime power conductors  $\mathfrak{f} := p^k$ . Thus, using subtractive sets as a challenge space for the verifier, one can argue that the extracted witness  $\bar{\mathbf{s}} := \bar{\mathbf{z}}/\bar{c}$  has low norm, because  $1/\bar{c}$  itself is short. However, this approach comes at a cost of non-negligible soundness error (due to the size of subtractive sets), and therefore some sort of soundness amplification is necessary. Furthermore, the protocol itself still does not manage to prove the exact norm bound, i.e.  $\|\mathbf{s}\| \leq \beta$ . In fact, in the context of recursive succinct arguments, the norm of the extracted witness can only be upper bounded by  $\gamma^\mu \cdot \theta^{O(\mu)} \cdot \beta$  for some  $\theta \approx \mathfrak{f}$ .

In the setting of power-of-two cyclotomic rings, the strategy above falls apart completely since there exists no subtractive set of size larger than two [Len76, AL21]. Hence, a different methodology has recently been developed. Notably, Beullens and Seiler [BS23] proposed a succinct argument, LaBRADOR, for proving  $\|\mathbf{s}\|^2 \leq \beta^2$  (among other relations), inspired by the following two-fold approach from [LNP22]:

- (i) *Approximate shortness proof.* Prove that  $\mathbf{s}$  is somewhat short.
- (ii)  $\mathbb{Z}_q$ -*Inner product proof.* Prove that  $(\langle \psi(\mathbf{s}), \psi(\mathbf{s}) \rangle \pmod q) \leq \beta^2$ , where  $\psi(\mathbf{s})$  is the coefficient vector of  $\mathbf{s}$ .

Combining (i) and (ii), one can argue that for a large enough modulus  $q$  no modulo wrap-around occurs, and therefore  $\langle \psi(\mathbf{s}), \psi(\mathbf{s}) \rangle \leq \beta^2$  holds over  $\mathbb{Z}$ .

In order to prove (i) without relying on subtractive sets, LaBRADOR uses the Johnson-Lindenstrauss random projection technique [BL17, LNS21, GHL22]. The idea is that the verifier will first generate an integer matrix  $\mathbf{B}$  with short (binary or ternary) values as a challenge, and the prover then outputs  $\psi(\mathbf{v}) := \mathbf{B}\psi(\mathbf{s}) \pmod q$ . Afterwards, the verifier checks whether  $\psi(\mathbf{v})$  is of low norm (which is true in the honest executions, since both  $\mathbf{B}$  and  $\psi(\mathbf{s})$  are). Finally, the prover needs to prove wellformedness of  $\psi(\mathbf{v})$ , i.e. the linear equation  $\mathbf{B}\psi(\mathbf{s}) = \psi(\mathbf{v})$  over  $\mathbb{Z}_q$ . The crucial soundness argument is that if the extracted  $\mathbf{s}$  was not short, then with high probability (dictated by the number of rows of  $\mathbf{B}$ ),  $\psi(\mathbf{v}) = \mathbf{B}\psi(\mathbf{s})$  would not have low norm, which leads to a contradiction. Unfortunately, the random projection strategy inherently requires the verifier to generate the matrix  $\mathbf{B}$ , which itself has length  $O(m)$ . As a consequence, the verifier

<sup>5</sup>For presentation, we omitted the factors related to the security parameter  $\lambda$ .

runtime becomes essentially linear in the witness size, which may not be satisfying in certain real-world use cases.

We highlight that both (i) and (ii) require some kind of inner product proof over  $\mathbb{Z}_q$ ; either between two committed vectors, or between one public and one committed vector. Since the underlying protocol natively operates over cyclotomic rings  $\mathcal{R} = \mathbb{Z}[\zeta]$ , it is essential to transform  $\mathbb{Z}$ -relations into equivalent ones over the ring  $\mathcal{R}$ . To this end, it was shown in [LNP22] that for any two elements  $a, b \in \mathcal{R}$  of a power-of-two cyclotomic ring, the constant term<sup>6</sup> of  $a \cdot \bar{b} \in \mathcal{R}$  is exactly equal to the inner product  $\langle \psi(\mathbf{a}), \psi(\mathbf{b}) \rangle \in \mathbb{Z}$ , where  $\psi(\mathbf{a}), \psi(\mathbf{b})$  are the coefficient vectors of  $a, b$  respectively and  $\bar{\cdot}$  here denotes the complex conjugation. This observation allows us to translate proving inner products and linear relations over integers into proving statements about constant terms over the ring  $\mathcal{R}$ . Finally, LaBRADOR makes use of the fact that inner product relations over  $\mathcal{R}$  are “folding-friendly” and can be efficiently proven in a recursive manner.

Interestingly, LaBRADOR also managed to circumvent the correctness gap by taking inspiration from the “decompose-then-hash” paradigm used in lattice-based Merkle trees [PSTY13]. Intuitively, using the notation above for describing recursive-based protocols, instead of folding the intermediate witness  $\mathbf{s}_{i-1}$  directly into a new one  $\mathbf{s}_i$ , an honest prover would first decompose  $\mathbf{s}_{i-1}$  (w.r.t. some decomposition base  $b$ ) into multiple vectors  $(\mathbf{s}_{i-1,j})_{j \in [\ell]}$  of much smaller norm and then fold all of them together into a new witness  $\mathbf{s}_i$ <sup>7</sup>. By carefully picking various parameters, such as  $b$ , one can ensure that, in an honest execution, if  $\|\mathbf{s}_{i-1}\| \leq \beta$ , then we must have  $\|\mathbf{s}_i\| \leq \beta$ . This technique was also adopted in a recent folding scheme called LatticeFold [BC24].

**Bridging the gap.** At a high level, the aforementioned approaches to prove shortness seem somewhat orthogonal. For  $\mathfrak{f} = p^k$ , where  $p = \text{poly}(\lambda)$  is a large enough prime, one can rely on subtractive sets to efficiently prove approximate shortness (i) with succinct verification [CLM23]. However, it is unknown how to translate proving  $\mathbb{Z}_q$ -relations, as in (ii), into equivalent relations over odd prime-power cyclotomic rings. On the other hand, for  $\mathfrak{f} = 2^k$ , one can apply the Johnson-Lindenstrauss projection strategy to prove both (i) and (ii), but at the cost of slow verification time.

Hence, it is an important research question whether there exist cyclotomic (or other) rings  $\mathcal{R}$ , which contain subtractive sets of fairly large size, and at the same time, expose efficient packing and batching techniques for turning relations over  $\mathbb{Z}$  (or more generally, other base rings) to relations over  $\mathcal{R}$ . An affirmative answer, together with existing optimisations, would then yield a practical lattice-based succinct argument for proving exact norm bounds with fast verification.

## 1.1 Our Contributions

In this work, we present a versatile toolkit for constructing lattice-based succinct arguments that eliminate correctness and soundness gaps while maintaining succinct verification. Our contributions are outlined as follows:

**Succinct Arguments for Bounded-Norm Satisfiability.** We design a lattice-based succinct argument system for bounded-norm satisfiability of structured linear and inner-product relations. Our system retains features of previous protocols, such as transparent setup, quasi-linear-time prover, and polylogarithmic-time verifier, while simultaneously eliminating any correctness and soundness gaps. Consequently, our argument system achieves asymptotically the most attractive proof sizes, which are smaller by at least a factor of  $\Omega(\log^2 \lambda)$  smaller than the prior state-of-the-art constructions (see Figure 1 for more details). Furthermore, our protocol’s modular design allows for straightforward analysis and customisation, making it adaptable to various applications.

**Subtractive Sets.** Our protocol uses subtractive sets as challenge sets. While subtractive sets for prime-power cyclotomic rings are well-known, the non-prime-power case seems less studied. Motivated by the need of non-prime-power rings (e.g. for the twisted trace technique, see below) in some applications, we identify a subtractive set for cyclotomic rings  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$  of non-prime-power conductor  $\mathfrak{f}$  with a cardinality of  $\mathfrak{f}/\mathfrak{f}_{\max}$ , where  $\mathfrak{f}_{\max}$  is the largest prime-power divisor of  $\mathfrak{f}$ . Additionally, we identify subtractive sets

<sup>6</sup>We say that  $a_0 \in \mathbb{Z}$  is the constant term of the ring element  $a = \sum_{i=0}^{\varphi(\mathfrak{f})} a_i \zeta^i \in \mathbb{Z}[\zeta]$ .

<sup>7</sup>For soundness, the prover needs to prove additional relations involving  $(\mathbf{s}_{i-1,j})_{j \in [\ell]}$ .

scheme	assumptions	proof size
[BCS23]	M-SIS	$O(\log^6 m \cdot \lambda^2 / \log \lambda)$
[CLM23]	vSIS	$O(\log^5 m \cdot \lambda^2 / \log^2 \lambda)$
[FMN23]	PowerBASIS	$O(\log^5 m \cdot \lambda^2 / \log^2 \lambda)$
[AFLN23]	M-SIS	$O(\log^5 m \cdot \lambda^2 / \log^2 \lambda)$
[CMNW24]	SIS	$O(\log^3 m \cdot \lambda^2)$
<b>This work</b>	vSIS	$O(\log^3 m \cdot \lambda^2 / \log^2 \lambda)$

**Fig. 1.** Asymptotic efficiency of our commitment opening proof (in bits) and comparison with prior interactive proofs which support succinct  $\text{poly}(\log m, \lambda)$  verification time. Here,  $\lambda$  is the security parameter and  $m$  is the length of the committed vector. For each construction, the proof size corresponds to the soundness error  $\text{poly}(\lambda, \log m) \cdot 2^{-\lambda}$ . The SIS-related parameters were chosen with respect to the methodology from [MR09] for running BKZ on block size  $b = O(\lambda)$ . For [BCS23, CLM23, FMN23] as well as our scheme, which only achieve inverse-polynomial soundness in one-shot, we applied a standard soundness amplification by parallel-repeating the protocol by a factor of  $O(\lambda / \log \lambda)$ . We highlight that for the sizes reported from [AFLN23, CMNW24], the knowledge extractor runs in expected super-polynomial time in  $m$  and  $\lambda$ .

over the real subrings  $\mathbb{Z}[\zeta_{\mathfrak{f}} + \zeta_{\mathfrak{f}}^{-1}]$ , with a cardinality of  $(p+1)/2$  for prime-power conductors  $\mathfrak{f} = p^k$  and  $\lfloor \mathfrak{f}/(2\mathfrak{f}_{\max}) \rfloor$  for non-prime-power  $\mathfrak{f}$ .

**Embedded  $\mathbb{Z}$ -Inner-Products via Twisted Trace.** While our protocol supports proving inner products over rings such as  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$ , higher-layer applications may require proving inner products over  $\mathbb{Z}$ , e.g. for proving that a committed  $\mathbb{Z}$ -vector is binary. Unfortunately, efficient methods for embedding  $\mathbb{Z}$ -inner products to  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$ -inner products were only known for  $\mathfrak{f} = 2^d$  being a power of 2, which is problematic because subtractive sets over  $\mathbb{Z}[\zeta_{2^d}]$  are of cardinality at most 2. We extend the existing embedding method to any ring of the form  $\mathbb{Z}[\zeta_{2^d}] \otimes \mathbb{Z}[\zeta_{p_0} + \zeta_{p_0}^{-1}] \otimes \dots \otimes \mathbb{Z}[\zeta_{p_{k-1}} + \zeta_{p_{k-1}}^{-1}]$ , where  $p_0, \dots, p_{k-1}$  are distinct odd primes. This is achieved by replacing the “constant term map” with a “twisted trace map” defined as:  $\tau(z) = \frac{1}{N} \text{Trace}(\alpha \cdot z)$ .

**Succinct Consistency Proof for CRT.** Another typical way of embedding  $\mathbb{Z}$ -relations into  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$ -relations is via the Chinese Remainder Transform (CRT). However, this requires proving that the witness vector is committed in both the coefficient embedding and its CRT coefficients consistently, and known consistency proofs are not succinct. Using the fact that the CRT over cyclotomic fields with smooth conductors can be succinctly represented through a few automorphism evaluations, we derive a succinct argument for the consistency between the commitment of the coefficient embedding and that of the CRT coefficients. At the core of our succinct consistency proof is a new succinct argument that verifies whether two committed vectors are related by an entry-wise automorphism.

## 2 Technical Overview

Throughout this work, we will assume that  $\mathcal{K} = \mathbb{Q}(\zeta)$  is a cyclotomic field with conductor  $\mathfrak{f}$  and degree  $\varphi = \varphi(\mathfrak{f}) = \text{poly}(\lambda)$ , and  $\mathcal{O}_{\mathcal{K}} = \mathbb{Z}[\zeta]$  is its ring of integers. For some of our results, we will further require  $\mathcal{K}^+ = \mathbb{Q}(\zeta + \zeta^{-1})$ , the maximal real subfield of  $\mathcal{K}$ , and its ring of integers  $\mathcal{O}_{\mathcal{K}^+} = \mathbb{Z}[\zeta + \zeta^{-1}]$ . Depending on the context of a specific section, we will use  $\mathcal{R} \subset \mathcal{O}_{\mathcal{K}}$  to denote a ring of interest to that section. Unless specified, we measure the norm of elements and vectors by their  $\ell_2$ -norm over the canonical embedding over  $\mathcal{K}$ . Our results can be divided into three parts, which we overview in Section 2.1, 2.2, and 2.3 respectively.

### 2.1 Subtractive Sets

In Section 4, we expose subtractive sets over  $\mathcal{O}_{\mathcal{K}}$  with non-prime-power conductor  $\mathfrak{f}$ , and over  $\mathcal{O}_{\mathcal{K}^+}$  with both prime-power and non-prime-power conductors, with favourable properties, i.e. they have  $\text{poly}(\lambda)$  cardinality and small expansion factors. These subtractive sets can be used in any lattice-based arguments, and in particular those developed in this work.

A set  $S \subset \mathcal{R}$  is said to be subtractive over  $\mathcal{R}$  if for any two distinct elements  $c, c' \in S$ , it holds that  $c - c' \in \mathcal{R}^\times$ , i.e.  $c - c'$  is a unit. This concept is prevalently linked with the examination of Euclidean number fields [Len76] and has also found relevance in lattice-based cryptography, specifically in argument systems and secret sharing [AL21]. An explicit creation of an upper-bound-matching cardinality  $p$  is evident in a cyclotomic ring  $\mathcal{R} = \mathcal{O}_\mathcal{K}$  with a prime-power conductor  $\mathfrak{f} = p^k$ . On the other hand, we are not aware of explicit studies of subtractive sets regarding other cyclotomic rings and their subrings.

For applications in lattice-based cryptography, the most relevant measures of the quality of a subtractive set  $S$  are its

- (i) cardinality  $|S|$ , which inversely affects the knowledge error of argument systems using  $S$  as a challenge set,
- (ii) “expansion factor”  $\gamma = \gamma_S$ , i.e. how much the norm of an element grows when multiplied with an element in  $S$ , which affects the “correctness gap” of lattice-based argument systems,
- (iii) “inverse-expansion factor”  $\theta = \theta_S$ , i.e. how much the norm of an element grow when multiplied with  $(c - c')^{-1}$  for distinct  $c, c' \in S$ , which affects the “soundness gap” of lattice-based argument systems.

For  $\mathcal{R} = \mathcal{O}_\mathcal{K}$  with prime-power conductor  $\mathfrak{f} = p^k$ , it is known [Len76, AL21] that there exists a subtractive set  $S$  of cardinality  $p$  and expansion factors  $\gamma, \theta \approx p$ .

Our main result in this part is the exposition of the subtractive set  $S := \{\zeta^i\}_{i \in [\mathfrak{f}/\mathfrak{f}_{\max}]}$  of cardinality  $|S| = \mathfrak{f}/\mathfrak{f}_{\max}$  for any conductor  $\mathfrak{f}$  with at least two distinct prime factors, where  $\mathfrak{f}_{\max}$  is the largest prime-power factor of  $\mathfrak{f}$ . Notably, the expansion factor is  $\gamma = 1$ , i.e. the norm of an element does not grow when multiplied with an element from  $S$ , while the inverse-expansion factor  $\theta \approx \mathfrak{f}$  is similar to the existing result for prime-power rings.

For completeness, we also expose related subtractive sets over  $\mathcal{O}_{\mathcal{K}+}$  for both prime-power and non-prime-power conductors.

## 2.2 Tight Succinct Argument for Bounded Norm Satisfiability

In Section 5, we work with  $\mathcal{R} = \mathcal{O}_\mathcal{K}$  or  $\mathcal{O}_{\mathcal{K}+}$ . We present a new lattice-based succinct argument for proving the bounded norm satisfiability of structured linear and/or inner-product relations, denoted by  $\Xi^{\text{lin}}$  and  $\Xi^{\text{ip}}$  respectively. More concretely, the argument system allows to prove knowledge of a short vector  $\mathbf{w} \in \mathcal{R}^m$ , with  $m = d^\mu$ , satisfying

- a linear relation  $\mathbf{F}\mathbf{w} = \mathbf{y} \bmod q$ , where  $\mathbf{F} = \mathbf{F}_{\mu-1} \bullet \dots \bullet \mathbf{F}_0 \in \mathcal{R}_q^{n \times m}$  can be expressed as a row-wise tensor product of  $\mu$  matrices  $\mathbf{F}_i \in \mathcal{R}_q^{n \times d}$ , and
- (optionally) an inner-product relation  $\langle \mathbf{w}, \alpha(\mathbf{w}) \rangle \bmod q$ , where  $\alpha$  is either the identity function or the complex conjugate (specified publicly).

Our argument system consists of  $O(\mu) = O(\log_d m)$  rounds and is public-coin, and can thus be made non-interactive via the Fiat-Shamir transform. The prover time is quasi-linear in the size of the statement, and both the proof size and the verifier time are polylogarithmic in the statement size. It can be instantiated with a transparent setup. For example, the rows of  $\mathbf{F}$  could contain a random commitment key of the vSIS commitment scheme [CLM23] and evaluations of monomials at different evaluation points. This turns the vSIS commitment scheme into a polynomial commitment scheme, which can then be used to compile a PIOP into a SNARK.

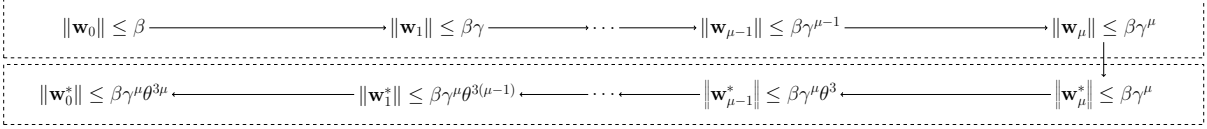
**Correctness and Soundness Gaps.** A distinguishing feature of our argument system is that it is free of the so-called “correctness gap” and “soundness gap”.

The correctness gap refers to the phenomenon that although the prover’s witness  $\mathbf{w}$  is of norm at most  $\beta$ , the norm check performed by the verifier in the protocol is against a bound  $\beta' \gg \beta$ . Typically, e.g. in lattice-based Bulletproofs, we have  $\beta' \approx (1 + \gamma)^\mu \beta$ . Using the subtractive set suggested in [AL21] and picking  $\mu \approx \log \lambda$ , the gap  $\beta'/\beta \approx (1 + \gamma)^\mu$  is super-polynomial in  $\lambda$ . Note that if the subtractive set suggested in Section 4 with  $\gamma = 1$  is used, then the correctness gap is immediately reduced to  $\text{poly}(\lambda)$  but still greater than 1 (i.e. no gap).

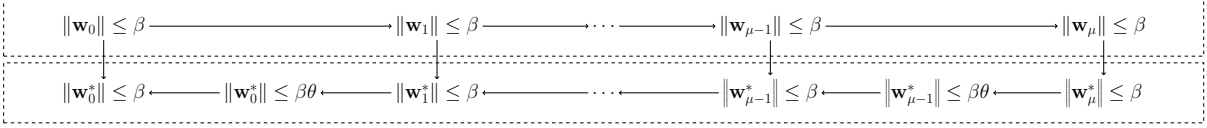
The more challenging issue is that of the soundness gap, which refers<sup>8</sup> to the limitation that, in addition to the correctness gap  $\beta'/\beta$ , the witness produced by a knowledge extractor is of even larger norm

<sup>8</sup>In general, the soundness gap consists of a “stretch”, i.e. increase in witness norm, and a “slack”, i.e. a multiplicative approximation factor. Using a subtractive set, the slack can be eliminated.

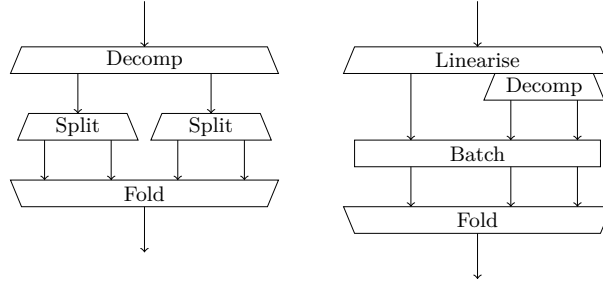
a) Lattice-based Bulletproofs.



b) Split-and-Fold + Norm-Check.



**Fig. 2.** Overview of the evolution of a prover witness  $\mathbf{w}_0$  to an extracted witness  $\mathbf{w}_0^*$  in lattice-based Bulletproofs and in Split-and-Fold + Norm-Check.



**Fig. 3.** Structures of Split-and-Fold (left) and Norm-Check (right) protocols.

$\beta^* \gg \beta'$ . Using the example of lattice-based Bulletproofs again, we have  $\beta^* \approx (2\theta)^{3\mu} \beta' \approx (1+\gamma)^\mu (2\theta)^{3\mu} \beta$ . Since no currently known subtractive set (including those suggested in Section 4) achieves  $\theta = O(1)$ , the soundness gap problem cannot be solved by simply using a different subtractive set, at least until more favourable sets are found.<sup>9</sup>

Figure 2 overviews the evolution of a prover witness  $\mathbf{w}_0$  to an extracted witness  $\mathbf{w}_0^*$  in lattice-based Bulletproofs and in this work.

**Lattice-based Bulletproofs.** In Fig. 2 part a) for Bulletproofs, each arrow in the top row represents one Bulletproofs folding step, where  $\mathbf{w}_i$  denotes the intermediate witness after the  $i$ -th folding step. The norm of the  $i$ -th round prover witness  $\mathbf{w}_i$  grows by a multiplicative factor of (around)  $\gamma$  compared to the previous round prover witness  $\mathbf{w}_{i-1}$ . The last round witness  $\mathbf{w}_\mu$  is then of norm around  $\beta\gamma^\mu$ , i.e. with correctness gap  $\gamma^\mu$ . The vertical arrow is trivial since the last-round prover witness is sent in plain, i.e.  $\mathbf{w}_\mu^* = \mathbf{w}_\mu$ . Each arrow in the bottom row represents a “traditional witness extraction step”, i.e. moving one layer up in the tree-special soundness witness extraction, where  $\mathbf{w}_i^*$  denotes the extracted witness at depth  $i$ . The norm of the  $i$ -th round extracted witness  $\mathbf{w}_i^*$  grows by (roughly) a multiplicative factor of  $\theta^3$  compared to the previous round extracted witness  $\mathbf{w}_{i-1}^*$ . The final extracted witness  $\mathbf{w}_0^*$  is then of norm around  $\beta\gamma^\mu\theta^{3\mu}$ , i.e. the soundness gap is  $\gamma^\mu\theta^{3\mu}$ .

**Split-and-Fold and Norm-Check.** To eliminate correctness and soundness gaps, we propose two modular protocols called Split-and-Fold and Norm-Check, each of which is a composition of atomic elementary building blocks – ( $b$ -ary) Decompose, Split, Fold, Linearise, and Batch. The structures of Split-and-Fold and Norm-Check, designed to eliminate the correctness and soundness gaps respectively, are depicted in Fig. 3. These protocols are designed to run in an interleaved manner to restrict the norm growth of the witness and to aid witness extraction. The lattice-based Bulletproofs protocol can be seen

<sup>9</sup>We believe that a slightly better but still super-polynomial soundness gap of  $\beta^*/\beta' \approx (1+\gamma)^\mu (2\theta)^\mu$  can be achieved using a technique called “short-circuit extraction” [HKR19].

as a repeated execution of the barebone Split protocol and Fold protocol without any norm-restricting mechanisms.

The Split-and-Fold and Norm-Check can be summarised as follows:

*Split-and-Fold.* The purpose of the Split-and-Fold protocol is to shrink the dimension of the relation to be proven without increasing the norm of the witness. On input a  $\Xi^{\text{lin}}$  instance  $(\mathbf{F}, \mathbf{y})$  with witness  $\mathbf{w}$  of norm at most  $\beta$ , run the Decompose protocol to decompose the witness into  $b$ -ary parts. This splits  $(\mathbf{F}, \mathbf{y})$  into  $\ell$  sub-instances  $(\mathbf{F}, \mathbf{y}_i)$  each with witness  $\mathbf{w}_i$ . If  $b$  is small, each sub-witness  $\mathbf{w}_i$  norms at most  $b/2 \ll \beta$ . Then, for each sub-instance  $(\mathbf{F}, \mathbf{y}_i)$  with witness  $\mathbf{w}_i$ , run the Split protocol to peel off one tensor factor of  $\mathbf{F}$ , i.e. factor  $\mathbf{F}$  into  $\mathbf{F} = \mathbf{R} \bullet \tilde{\mathbf{F}}$  and decompose  $\mathbf{w}_i$  into  $(\mathbf{w}_{i,j})_{j \in [d]}$ . Each sub-instance is thus further split into finer sub-instances  $(\tilde{\mathbf{F}}, \mathbf{y}_{i,j})$  for some appropriate  $\mathbf{y}_{i,j}$  with witnesses  $\mathbf{w}_{i,j}$  still of norm at most  $b/2$ . Finally, the Fold protocol is run to merge all sub-instances into a single instance  $(\tilde{\mathbf{F}}, \tilde{\mathbf{y}})$  with witness  $\tilde{\mathbf{w}}$ . For appropriately chosen  $b$ , we should end up with a next-round witness  $\tilde{\mathbf{w}}$  of norm at most  $\beta$  again.

Note that the above suffices to eliminate the correctness gap, i.e. if the Split-and-Fold protocol were to be run recursively, the intermediate and hence final witnesses of the prover will remain to have norm at most  $\beta$ . However, given an extracted next-round witness  $\tilde{\mathbf{w}}$  of norm at most  $\beta$ , the knowledge extractor could only extract a candidate witness of norm at most  $\approx \theta\beta$ . To improve the norm bound to  $\beta$ , we need to interleave a Norm-Check protocol, as described below, between executions of Split-and-Fold.

*Norm-Check.* The purpose of the Norm-Check protocol is to upgrade a relaxed norm bound guarantee to a tight norm bound guarantee. On input a  $\Xi^{\text{lin}}$  instance  $(\mathbf{F}, \mathbf{y})$  with witness  $\mathbf{w}$  of norm at most  $\beta$ , the prover sends the value  $t = \langle \mathbf{w}, \bar{\mathbf{w}} \rangle_{\mathcal{R}}$  to the verifier, who can check that  $\text{Trace}(t) \leq \beta^2$ . If  $t$  is computed correctly, then  $\text{Trace}(t)$  is precisely the square of the canonical  $\ell_2$ -norm of  $\mathbf{w}$ , which the verifier checks to be at most  $\beta^2$ . It thus remains for the prover to prove that  $t$  is computed correctly, along with all the other relations. To do so, the prover encodes  $\mathbf{w}$  as the coefficients of a polynomial  $g(X)$ , and commit to the coefficients of the Laurent polynomial  $L(X) = g(X) \cdot \bar{g}(X^{-1})$ . This reduces the problem to checking that  $L$  is computed correctly and has constant term  $t$ , both of which can be expressed as relations captured by  $\Xi^{\text{lin}}$ .

Two issues remain: First, the norm of the coefficients of  $L(X)$  is around  $\beta^2$  instead of  $\beta$ . To tackle this, the prover runs the Decompose protocol (in a non-black-box manner) to shrink the coefficients of  $L(X)$  back to norm  $\beta$ , at the cost of spawning new sub-instances. Second, the extra checks for  $L$  being computed correctly and  $L(0) = t$  introduce more constraints which would translate to higher communication cost when handled naively. To tackle this, the parties run the Batch protocol to compress the newly added constraints with the existing ones. Finally, we use again the Fold protocol to merge all sub-instances into one.

In Fig. 2 part b) for “Split-and-Fold + Norm-Check”, each horizontal arrow in the top row represents one “split-and-fold” step which replaces a Bulletproofs folding step. The effect of this is that the norm of  $\mathbf{w}_i$  remains at most  $\beta$  for all  $i$ , i.e. the correctness gap is eliminated. Each vertical arrow from  $\mathbf{w}_i$  to  $\mathbf{w}_i^*$  for  $i < \mu$  represents one “Norm-Check” which is used to prove that the norm of the current witness  $\mathbf{w}_i$  is at most  $\beta$ . For this step, it is important that the norm function is chosen to be the  $\ell_2$ -norm over the canonical embedding, so that it can be expressed in terms of the (complex) inner product which is natively supported by our protocol. This proof upgrades the bound  $\|\mathbf{w}_i^*\| \leq \theta\beta$  guaranteed by a traditional witness extraction step to a tighter bound  $\|\mathbf{w}_i^*\| \leq \beta$ , i.e. the soundness gap is eliminated.

### 2.3 Embedding $\mathbb{Z}$ -Inner Products

Lattice-based succinct arguments such as those constructed in Section 5 typically support proving relations over a ring  $\mathcal{R}$  natively. However, in many applications, we would like to prove algebraic statements given over  $\mathbb{Z}$ , which motivates the question of how to reduce a statement over  $\mathbb{Z}$  to statements over  $\mathcal{R}$ , so that a proof of the latter implies a proof of the former. Specifically, we consider the task of proving that some (committed) vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{Z}^{m^\delta}$  satisfies  $\langle \mathbf{x}, \mathbf{y} \rangle = z$  for some given  $z \in \mathbb{Z}$ . This task is of particular interest since, for some applications (e.g. constructing verifiable delay function [LM23]) it is necessary for the prover to prove that the witness is not only short but in fact binary. More generally, the application might require the prover to show a proof for  $\mathbf{x} \in [a, b]^{m^\delta}$  for some  $a, b \in \mathbb{Z}$ , which is not immediately implied by a bounded-norm guarantee.

To prove binariness, the basic idea is, for a witness  $\mathbf{w} \in \mathbb{Z}^{m^\delta}$ , to use the equivalence  $\mathbf{w} \in \{0, 1\}^{m^\delta} \iff \langle 1^{m^\delta} - \mathbf{w}, \mathbf{w} \rangle_{\mathbb{Z}} = 0$  to reduce checking the binariness of  $\mathbf{w}$  to checking that some transformed witness

vector over  $\mathcal{R}$  is short and satisfies some linear and inner-product relations, where  $\mathcal{R} \subset \mathcal{O}_{\mathcal{K}}$  is of dimension  $\delta \mid \varphi$  when viewed as  $\mathbb{Z}$ -modules.

**Existing Embedding Methods.** We are aware of three ways to embed  $\mathbb{Z}$ -inner products into  $\mathcal{R}$ -inner products in the literature, each with a significant drawback:

- (i) Naive embedding: Interpret each  $\mathbb{Z}$  element as an  $\mathcal{R}$  element via the inclusion  $\mathbb{Z} \subset \mathcal{R}$ , and interpret the  $\mathbb{Z}$ -inner product as an  $\mathcal{R}$ -inner product. This incurs a multiplicative overhead of  $\delta$  in terms of statement and witness sizes, which translate into overheads in prover and verifier computation, proof size, etc.
- (ii) Coefficient embedding: Divide the witness into blocks containing  $\delta$   $\mathbb{Z}$ -elements, and encode each block as an  $\mathcal{R}$  element via the (inverse-)coefficient embedding<sup>10</sup>  $\psi^{-1} : \mathbb{Z}^{m\delta} \rightarrow \mathcal{R}^m$ . For certain  $\mathcal{R}$ , we have

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \text{ct}(\langle \psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{y})} \rangle_{\mathcal{R}})$$

where  $\text{ct}(\cdot)$  denotes the constant term of the coefficient embedding.

This embedding has a convenient property that it is (somewhat) norm-preserving, i.e.  $\mathbf{x}$  is short if and only if  $\psi^{-1}(\mathbf{x})$  is also short (in both coefficient and canonical embedding). However, this approach only works for  $\mathbb{Z}[\zeta_{2^d}]$ . This is problematic since the largest subtractive set over  $\mathbb{Z}[\zeta_{2^d}]$  is  $\{0, 1\}$ .

- (iii) CRT embedding: Let the witness vectors be such that  $\mathbf{x}, \mathbf{y} \in \mathbb{Z}_p^{m\delta}$  for some (typically small) prime  $p$  which splits completely in  $\mathcal{R}$ . Divide the witness into blocks of  $\delta$   $\mathbb{Z}$  elements, and encode each block as an  $\mathcal{R}$  element via the (inverse-)CRT embedding  $\text{CRT}_p^{-1} : \mathbb{Z}_p^{m\delta} \rightarrow \mathcal{R}_p^m$ . It holds that

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \langle \mathbf{1}^{\delta}, \text{CRT}_p(\langle \text{CRT}_p^{-1}(\mathbf{x}), \text{CRT}_p^{-1}(\mathbf{y}) \rangle_{\mathcal{R}}) \rangle_{\mathbb{Z}} \bmod p.$$

This approach is powerful in that it not only supports proving about  $\mathbb{Z}_p$ -inner products, but in fact about  $\mathbb{Z}_p$ -Hadamard products  $\mathbf{x} \odot \mathbf{y} \bmod p$ , which is more fine-grained. However, to turn a claim about  $\mathbb{Z}_p$ -inner products into a claim about  $\mathbb{Z}$ -inner products (without reduction modulo  $p$ ), we would additionally need to prove that  $\|\langle \mathbf{x}, \mathbf{y} \rangle\|_{\infty} < p/2$ , so that the reduction modulo  $p$  has no effect. Since  $\text{CRT}_p$  does not respect the geometry of  $\mathbb{Z}$  and  $\mathcal{R}$ , this approach usually requires the prover to commit to the witness vectors in both the  $\psi^{-1}(\cdot)$  and  $\text{CRT}_p^{-1}(\cdot)$  encodings, prove that the former is short, and prove that the two commitments are consistent. An issue here is that existing proofs of consistency between the two encodings (e.g. [BS23, LNS20]) do not have a succinct verifier, i.e. they run in time linear in the witness size.

In the following, we highlight how the aforementioned issues regarding the coefficient and CRT embeddings can be solved over certain (wide) range of rings.

**Twisted Trace Maps.** In Section 6, we generalise the coefficient embedding technique over power-of-2 rings to a wide range of other rings. Recall from the above that, over  $\mathcal{O}_{\mathcal{K}}$  with a power-of-2 conductor, it holds that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \text{ct}(\langle \psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{y})} \rangle_{\mathcal{R}})$ . In fact, the constant term function can be expressed as  $\text{ct}(\cdot) = \frac{1}{\varphi} \cdot \text{Trace}_{\mathcal{K}/\mathbb{Q}}(\cdot)$  where  $\text{Trace}_{\mathcal{K}/\mathbb{Q}}$  denotes the field trace, and the power basis  $\{1, \zeta, \dots, \zeta^{\varphi-1}\}$  satisfies i.e. the power basis is orthogonal with respect to the field trace.

The above point of view motivates the search for ideal lattices with  $\mathbb{Z}$ -bases orthogonal with respect to the field trace. This leads us to the literature of lattice constellations. In particular, we extract the following embedding method from [BFOV04]: Over  $\mathcal{O}_{\mathcal{K}^+}$  with prime conductor  $\mathfrak{f}$ , there exists an (efficiently computable) basis  $\mathbf{b}^+ \in \mathcal{O}_{\mathcal{K}^+}^{\varphi/2}$  and a twist element  $\alpha \in \mathcal{O}_{\mathcal{K}^+}$  such that

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \frac{1}{2\mathfrak{f}} \text{Trace}_{\mathcal{K}/\mathbb{Q}}(\alpha \cdot \langle \psi_{\mathbf{b}^+}^{-1}(\mathbf{x}), \overline{\psi_{\mathbf{b}^+}^{-1}(\mathbf{y})} \rangle_{\mathcal{R}})$$

where  $\psi_{\mathbf{b}^+} : \mathcal{O}_{\mathcal{K}^+} \rightarrow \mathbb{Z}^{\varphi/2}$  denotes the coefficient embedding with respect to the basis  $\mathbf{b}^+$ . Furthermore, adapting a result from the same work [BFOV04] regarding tensor products of rings, we extract similar embedding methods based on twisted trace maps for rings  $\mathcal{R}$  of the form  $\mathcal{R} = \mathcal{O}_{\mathcal{K}_{2^d}} \otimes \mathcal{O}_{\mathcal{K}_{p_0}^+} \otimes \dots \otimes \mathcal{O}_{\mathcal{K}_{p_{k-1}}^+}$ , where the subscripts of  $\mathcal{K}$  denote the conductors the respective factor rings and  $p_0, \dots, p_{k-1}$  are distinct odd primes. This captures power-of-2 rings as a special case. Notably, since such  $\mathcal{R}$  generally have non-prime-power conductors, they are compatible with the subtractive set for non-prime-power rings exposed in Section 4.

<sup>10</sup>For example, with respect to the power basis  $\{1, \zeta, \dots, \zeta^{\varphi-1}\}$  of a cyclotomic field, the coefficient embedding of an element  $x = \sum_{i \in [\varphi]} x_i \zeta^i$  is denoted as  $\psi(x) = (x_i)_{i \in [\varphi]}$ .



**Succinct Proof for Consistency of CRT.** As highlighted earlier, the missing piece, required to harness the power of the CRT embedding for Hadamard and inner products, is a verifier-succinct argument for proving the consistency between the coefficient embedding and the CRT embedding. More precisely, we need a succinct argument for proving that two ring vectors  $\mathbf{w}, \mathbf{w}' \in \mathcal{R}^m$  satisfy

$$\psi(\mathbf{w}) = \text{CRT}_p(\mathbf{w}') \bmod p. \quad (1)$$

In Section 7, we present a protocol for performing this task over  $\mathcal{R} = \mathcal{O}_{\mathcal{K}}$  where the conductor  $\mathfrak{f}$  is  $w$ -smooth, i.e. all its prime factors are at most some small integer  $w$ , with proof size and verifier time scaling linearly in  $w \log_w \mathfrak{f}$ . In other words, if  $w = O(1)$ , then the complexity is logarithmic in  $\mathfrak{f}$ .

Underlying our protocol is the observation that, if the conductor  $\mathfrak{f}$  is  $w$ -smooth, then the map  $\text{CRT}_p^{-1} \circ \psi$  can be expressed as the composition of  $t \leq O(\log \mathfrak{f})$  maps, each being a linear combination of  $h \leq O(\log \mathfrak{f})$  automorphisms from  $\text{Gal}(\mathcal{K}/\mathbb{Q})$  with coefficients lying in  $\mathcal{R}$ . This means that, to succinctly prove that  $\mathbf{w}' = \text{CRT}_p^{-1}(\psi(\mathbf{w})) \bmod p$ , it suffices to design a succinct argument for proving automorphism relations.

Motivated by the above, we present a succinct reduction of knowledge from checking the automorphism relation  $\alpha(\mathbf{w}) = \mathbf{w}'$  to checking that  $(\mathbf{w}, \mathbf{w}')$  satisfies some linear relations. Combined with the Split-and-Fold and Norm-Check protocols designed in Section 5, we obtain a succinct argument for proving Eq. (1).

### 3 Preliminaries

Let  $\mathbb{N} = \{1, 2, \dots\}$  denotes natural numbers and  $\lambda \in \mathbb{N}$  be the security parameter. For  $n \in \mathbb{N}$ , we write  $[n] := \{0, \dots, n-1\}$  counting from 0. For multidimensional ranges, we use the shorthand  $(i, j, k) \in [n, m, \ell]$  for  $i \in [n]$ ,  $j \in [m]$ , and  $k \in [\ell]$ .

Throughout this work, we let  $\mathcal{K} = \mathbb{Q}(\zeta)$  be a cyclotomic field with conductor  $\mathfrak{f}$  of degree  $\varphi = \varphi(\mathfrak{f})$ , where  $\zeta$  is a root of unity of order  $\mathfrak{f}$  and  $\varphi$  is Euler's totient function, and  $\mathcal{O}_{\mathcal{K}} = \mathbb{Z}[\zeta]$  be its ring of integers. We will also consider the maximal real subfield  $\mathcal{K}^+ = \mathbb{Q}(\zeta + \zeta^{-1})$  of  $\mathcal{K}$  and its ring of integers  $\mathcal{O}_{\mathcal{K}^+} = \mathbb{Z}[\zeta + \zeta^{-1}]$ . In contexts where we refer to multiple cyclotomic fields with different conductors  $(\mathfrak{f}_i)_{i \in [k]}$ , we write  $\mathcal{K}_{\mathfrak{f}_i}$  for  $i \in [k]$  to emphasise the conductors. We will usually use  $\mathcal{R} \subseteq \mathcal{O}_{\mathcal{K}}$  to denote a subring which has dimension  $\delta$  when viewed as a  $\mathbb{Z}$ -module. In Section A.1, we recall some basics of algebraic number theory.

#### 3.1 Cryptographic Assumption

We state an equivalent formulation of the vanishing short integer solution (vSIS) assumption [CLM23], which has a simpler description and better aligns with the notation adopted in this work. For more discussion on vSIS, we refer to Section A.2.

**Definition 1 (vSIS Assumption (adapted from [CLM23])).** *Let  $\text{params} = (\mathcal{R}, q, \beta, \chi)$  be parametrised by  $\lambda$ , where  $\mathcal{R}$  is a ring,  $q \in \mathbb{N}$  a modulus,  $\beta > 0$  a norm bound, and  $\chi$  a distribution over  $\mathcal{R}_q^{n \times \otimes_{i \in [\mu]} d_i}$  for some dimensions  $n, d_0, \dots, d_{\mu-1}, \mu \in \mathbb{N}$ . The  $\text{vSIS}_{\text{params}}$  assumption states that, for any PPT adversary  $\mathcal{A}$ , the advantage function satisfies*

$$\text{Adv}_{\text{params}, \mathcal{A}}^{\text{vSIS}}(\lambda) := \Pr \left[ \begin{array}{l} \mathbf{F}\mathbf{w} = \mathbf{0} \bmod q \\ \|\sigma(\mathbf{w})\|_2 \leq \beta \end{array} \middle| \begin{array}{l} \mathbf{F} \leftarrow_{\$} \chi \\ \mathbf{w} \leftarrow \mathcal{A}(\mathbf{F}) \end{array} \right] \leq \text{negl}(\lambda).$$

For simplicity, in this work, we will consider the setting where the block sizes  $d_0, \dots, d_{\mu-1}$  are identically set to some  $d \in \mathbb{N}$ , so that  $\mathbf{F}$  can be factored into  $\mathbf{F} = \mathbf{F}_{\mu-1} \bullet \dots \bullet \mathbf{F}_0$  with  $\mathbf{F}_i \in \mathcal{R}_q^{n \times d}$ , where  $\bullet$  denotes the row-wise tensor product.

#### 3.2 Reduction of Knowledge

In this paper we consider ternary relations  $\Xi \subseteq \{0, 1\}^* \times \{0, 1\}^* \times \{0, 1\}^*$ , where a tuple  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi$  consists of public parameters  $\text{pp}$ , statement  $\text{stmt}$  and witness  $\text{wit}$ . For presentation, we omit including  $\text{pp}$  when it is known from the context. We consider a modified and simplified definition of a reduction of knowledge [KP23] for the following reasons: All of our protocols are *public coin* and (*coordinate-wise*)

special sound [FMN23] or similar.<sup>11</sup> Thus, public reducibility is automatic and we have (super-constant) sequential composition results due to known (tree) black-box extractors, whereas composition in [KP23] is limited a constant number of protocols. Lastly, we define a *relaxed* knowledge soundness notion which is not present in [KP23]. For lack of space, we provide a condensed overview of reductions of knowledge. See Section A.3 for details.

**Definition 2 (Reduction of Knowledge (modified)).** *Let  $\Xi_0, \Xi_1$  be ternary relations. A reduction of knowledge (RoK)  $\Pi$  from  $\Xi_0$  to  $\Xi_1$ , short  $\Pi: \Xi_0 \rightarrow \Xi_1$ , is defined by two PPT algorithms  $\Pi = (\mathcal{P}, \mathcal{V})$ , the prover  $\mathcal{P}$ , and the verifier  $\mathcal{V}$ , with the following interface:*

- $\mathcal{P}(\text{pp}, \text{stmt}_1, \text{wit}_1) \rightarrow (\text{stmt}_2, \text{wit}_2)$ : *Interactively reduce the input statement  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi_0$  to a new statement  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi_1$  or  $\perp$ .*
- $\mathcal{V}(\text{pp}, \text{stmt}) \rightarrow \text{stmt}$ : *Interactively reduce the task of checking the input statement  $(\text{pp}, \text{stmt})$  w.r.t  $\Xi_0$  to checking a new statement  $(\text{pp}, \text{stmt})$  w.r.t.  $\Xi_1$ .*

A RoK  $\Pi$  is *correct*, if for any honest protocol run (with correct inputs), the prover outputs a witness for the reduced statement (which the verifier outputs). A RoK  $\Pi$  is *relaxed knowledge sound* from  $\Xi_0^{\text{KS}}$  to  $\Xi_1^{\text{KS}}$  with knowledge error  $\kappa(\text{pp}, \text{stmt})$  if there is a *black-box* expected polynomial-time extractor  $\mathcal{E}$ , which succeeds with probability  $\epsilon - \kappa(\text{pp}, \text{stmt})$  if the malicious prover outputs a valid witness for the reduced statement with probability  $\epsilon$  (on verifier's input  $(\text{pp}, \text{stmt})$ ).

## 4 Subtractive Sets

A subtractive set  $S$  over a ring  $\mathcal{R}$  is such that  $c - c'$  is a unit for any distinct  $c, c' \in S$ . While the notion is connected to the study of Euclidean number fields [Len76], it also found applications in lattice-based cryptography in the contexts of argument systems and secret sharing [AL21]. For a cyclotomic ring  $\mathcal{R}$  with prime-power conductor  $\mathfrak{f} = p^k$ , an explicit construction of upper-bound-matching cardinality  $p$  is known. For other cyclotomic rings and their subrings, however, not much seem to be explicitly studied. In this section, we construct subtractive sets over non-prime-power cyclotomic rings, as well as *real* cyclotomic rings.

**Definition 3 (Subtractive Set).** *We say that a set  $S \subseteq \mathcal{R}$  is subtractive over  $\mathcal{R}$  if  $c - c' \in \mathcal{R}^\times$  for any distinct  $c, c' \in S$ .*

While [AL21] measured the quality of a subtractive set over cyclotomic rings in terms of the  $\ell_\infty$ -norm over the coefficient embedding, in this work, we will instead work with the  $\ell_\infty$ -norm over the canonical embedding for compatibility with Section 5 via the inequality  $\forall c, x \in \mathcal{R}, \|\sigma(c \cdot x)\|_2 \leq \|\sigma(c)\|_\infty \cdot \|\sigma(x)\|_2$ . We measure the quality of a subtractive set by its cardinality, expansion factor  $\gamma_S$ , and inverse-expansion factor  $\theta_S$ , with the latter two defined below.

**Definition 4 ((Inverse-)Expansion Factor of Subtractive Set).** *Let  $S \subseteq \mathcal{R}$  be subtractive over  $\mathcal{R}$ . The expansion and inverse-expansion factors of  $S$  are*

$$\gamma_S := \max_{c \in S} \|\sigma(c)\|_\infty \quad \text{and} \quad \theta_S := \max_{c, c' \in S, c \neq c'} \left\| \sigma\left(\frac{1}{c - c'}\right) \right\|_\infty$$

respectively.

The following lemma often is handy for analysing inverse-expansion factors.

**Lemma 1.** *Let  $K = \mathbb{Q}(\zeta)$  with  $\zeta$  a primitive  $\mathfrak{f}$ -th root of unity such that  $\mathfrak{f} > 4$ . It holds that  $\left\| \sigma\left(\frac{1}{1-\zeta}\right) \right\|_\infty \leq \frac{\mathfrak{f}}{4\sqrt{2}}$ . Furthermore, if  $\zeta^i$  is a  $k$ -th root of unity, then  $\left\| \sigma\left(\frac{1}{1-\zeta^i}\right) \right\|_\infty \leq \frac{k}{4\sqrt{2}}$  and  $\left\| \sigma\left(\frac{1}{1+\zeta^i}\right) \right\|_\infty \leq \frac{k}{2\sqrt{2}}$ .*

<sup>11</sup> To turn soundness errors of probabilistic tests (such as Schwartz-Zippel) into knowledge errors, we merely need uniformly random transcripts. These are produced by (CW)SS extractors for example. We call such extractors *k-transcript extractors*.

*Proof.* By the definition of  $\|\sigma(\cdot)\|_\infty$ , we need to upper bound  $\max_{\sigma_j} \left| \sigma_j \left( \frac{1}{1-\zeta} \right) \right| = \max_{\sigma_j} \left| \frac{1}{1-\sigma_j(\zeta)} \right|$ , where  $\sigma_j$  ranges from  $\text{Gal}(\mathcal{K}/\mathbb{Q})$ . Since  $\sigma_j(\zeta)$  ranges over all primitive  $\mathfrak{f}$ -th root of unity, this is the same as  $\max_{j \in \mathbb{Z}_\mathfrak{f}^\times} \left| \frac{1}{1-\zeta^j} \right| = \max_{j \in \mathbb{Z}_\mathfrak{f}^\times} \left| \frac{1}{1-e^{j \cdot 2\pi i/\mathfrak{f}}} \right|$ . Thus, it suffices to lower-bound  $|1 - e^{j \cdot 2\pi i/\mathfrak{f}}|$  over  $j \in \mathbb{Z}_\mathfrak{f}^\times$ . Geometrically,  $e^{j \cdot 2\pi i/\mathfrak{f}}$  are points on the unit circle in the complex plane with angles incremented by  $2\pi j/\mathfrak{f}$ . Thus, the value is approximately  $|1 - e^{j \cdot 2\pi i/\mathfrak{f}}| \approx 2\pi j/\mathfrak{f}$  for small  $2\pi j/\mathfrak{f}$ . For an explicit bound, observe that for  $\alpha \leq \frac{1}{4}$  we have  $|1 - e^{\alpha 2\pi i}| = 2 \cdot \sin(\frac{2\pi\alpha}{2}) \geq \alpha \cdot 4\sqrt{2}$ . Setting  $\alpha = \frac{1}{\mathfrak{f}}$  proves the claim. Observe that the above argument only depends on the multiplicative order of  $\zeta$ , thus, the claim  $\left\| \sigma \left( \frac{1}{1-\zeta^i} \right) \right\|_\infty \leq \frac{k}{4\sqrt{2}}$  follows for  $\zeta^i$  being a  $k$ -th root of unity.

Next, we show that  $\left\| \sigma \left( \frac{1}{1+\zeta^i} \right) \right\|_\infty \leq \frac{k}{2\sqrt{2}}$  for  $\zeta^i$  being a  $k$ -th root of unity. Observe that either  $-\zeta^i$  is already a  $k$ -th root unity, then the proof concludes by applying the above bound for  $\left\| \sigma \left( \frac{1}{1-(-\zeta^i)} \right) \right\|_\infty$ . Otherwise,  $-\zeta^i$  is a  $2k$ -th root unity. Thus  $\left\| \sigma \left( \frac{1}{1+\zeta^i} \right) \right\|_\infty = \left\| \sigma \left( \frac{1}{1-(-\zeta^i)} \right) \right\|_\infty \leq \frac{2k}{4\sqrt{2}} = \frac{k}{2\sqrt{2}}$ .  $\square$

#### 4.1 Prime-Power Cyclotomics

We recall the subtractive set for prime-power cyclotomics [Len76, AL21] with conductor  $\mathfrak{f} = p^k$  and analyse its (inverse-)expansion factor in canonical  $\ell_2$ -norm. Although we are interested mostly in  $p \gg 2$ , the result also holds for  $p = 2$ .

**Theorem 1.** *Let  $\mathfrak{f} = p^k > 4$  for some prime  $p$ . The set  $S := \{\mu_0, \dots, \mu_{p-1}\} \subseteq_p \mathcal{O}_\mathcal{K}$  is subtractive, where  $\mu_i = (\zeta^i - 1)/(\zeta - 1)$ . Further,  $\gamma_S \leq p$  and  $\theta_S \leq \frac{\mathfrak{f}}{2\sqrt{2}}$ .*

*Proof.* Let  $i < j \in [p]$ . Observe that  $\mu_j - \mu_i = \zeta^i + \zeta^{i+1} + \dots + \zeta^{j-1} = \zeta^i \cdot \frac{\zeta^{j-i} - 1}{\zeta - 1}$  which is clearly a unit in  $\mathcal{R}$ , hence  $S$  is subtractive.

Below, write  $\|\cdot\|$  for  $\|\sigma(\cdot)\|_\infty$ . For the expansion factor, note that  $\mu_i$  is a sum of  $i$  roots of unity and  $i < p$ . Therefore  $\gamma_S = \max_{i \in [p]} \|\mu_i\| < p$ . For the inverse-expansion factor, observe that

$$\left\| \frac{1}{\mu_j - \mu_i} \right\| = \left\| \zeta^{-i} \cdot \frac{\zeta - 1}{\zeta^{j-i} - 1} \right\| \leq \|\zeta - 1\| \cdot \left\| \frac{1}{\zeta^{j-i} - 1} \right\| \leq 2 \left\| \frac{1}{\zeta^{j-i} - 1} \right\| \leq \frac{\mathfrak{f}}{2\sqrt{2}},$$

where the last inequality follows from Lemma 1 and the rest are elementary.  $\square$

#### 4.2 Non-Prime-Power Cyclotomics

A drawback of the subtractive set recalled above is its rather large expansion factor  $\gamma_S \leq p$ . In some applications, e.g. Section 5, we would like  $\gamma_S$  to be constant. Below, we expose a subtractive set over non-prime-power cyclotomic rings with expansion factor  $\gamma_S = 1$ .

**Theorem 2.** *Let  $\mathfrak{f}$  factor into  $k \geq 2$  coprime prime-power factors  $(\hat{\mathfrak{f}}_i)_{i \in [k]}$ , i.e.  $\mathfrak{f} = \prod_{i \in [k]} \hat{\mathfrak{f}}_i$ . Write  $\hat{\mathfrak{f}}_{\max} := \max_{i \in [k]} \hat{\mathfrak{f}}_i$ . The set  $S := \left\{ 1, \zeta, \zeta^2, \dots, \zeta^{\hat{\mathfrak{f}}_{\max}-1} \right\} \subseteq_{\mathfrak{f}/\hat{\mathfrak{f}}_{\max}} \mathcal{O}_\mathcal{K}$ , is subtractive. Furthermore,  $\gamma_S = 1$  and  $\theta_S \leq \frac{\mathfrak{f}}{4\sqrt{2}}$ .*

To prove Theorem 2, we begin with the following lemma which we believe should be well-established together with a supportive proposition. Since we could not find an explicit reference to the lemma, we provide a proof.

**Lemma 2.** *Let  $\mathcal{R} = \mathbb{Z}[\zeta_\mathfrak{f}]$  with a conductor  $\mathfrak{f}$  having  $k \geq 2$  coprime prime-power factors<sup>12</sup>  $(\hat{\mathfrak{f}}_i)_{i \in [k]}$ , i.e.  $\mathfrak{f} = \prod_{i \in [k]} \hat{\mathfrak{f}}_i$ . Write  $\hat{\mathfrak{f}}_{\max} := \max_{i \in [k]} \hat{\mathfrak{f}}_i$ . For  $j \in \left\{ 1, 2, \dots, \frac{\mathfrak{f}}{\hat{\mathfrak{f}}_{\max}} - 1 \right\}$ , it holds that  $1 - \zeta^j \in \mathcal{R}^\times$ .*

<sup>12</sup>For example,  $(2^3, 3^2)$  are coprime prime-power factors of  $72 = 2^3 3^2$ , but  $(2, 2^2, 3^2)$  are not.

*Proof.* Write  $\zeta = \zeta_{\mathfrak{f}}$ . First, consider the case when  $\zeta^j$  is a primitive  $\mathfrak{f}$ -th root of unity. Then, by Proposition 1,  $1 - \zeta^j$  is a unit in  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$ . If  $\zeta^j$  is not a primitive  $\mathfrak{f}$ -th root of unity, then it is a primitive  $\mathfrak{h}$ -th root of unity for some  $\mathfrak{h} \mid \mathfrak{f}$  and  $\zeta^j \in \mathbb{Z}[\zeta_{\mathfrak{h}}]$ . Observe that  $\frac{\mathfrak{f}}{\mathfrak{h}} \mid j$ . Assume that  $\mathfrak{h}$  is a prime-power, i.e.  $\mathfrak{h} = \hat{\mathfrak{f}}_i^n$  for some  $i \in [k]$  and  $n \geq 2$ . Hence, as  $j \in \left\{1, 2, \dots, \frac{\mathfrak{f}}{\hat{\mathfrak{f}}_{\max}} - 1\right\}$ ,

$$\frac{\mathfrak{f}}{\hat{\mathfrak{f}}_i^n} \leq j < \frac{\mathfrak{f}}{\hat{\mathfrak{f}}_{\max}},$$

which implies  $\hat{\mathfrak{f}}_{\max} < \hat{\mathfrak{f}}_i^n$ , a contradiction. Therefore,  $\mathfrak{h}$  is not a prime power, i.e. it has more than one distinct prime factors. By Proposition 1,  $1 - \zeta^j$  is invertible in  $\mathcal{R}_{\mathfrak{h}}$ , thus in  $\mathcal{R}_{\mathfrak{f}}$ .  $\square$

Next, we recall an elementary result.

**Proposition 1 ([Was97, Proposition 2.8]).** *Suppose  $\mathfrak{f}$  has at least two distinct prime factors. Then,  $1 - \zeta$  is a unit in  $\mathcal{R} = \mathbb{Z}[\zeta_{\mathfrak{f}}]$  for any  $\mathfrak{f}$ -th primitive root of unity  $\zeta$ .*

Finally, we state our proof of Theorem 2.

*Proof.* (Theorem 2) For  $i, j \in [\mathfrak{f}/\hat{\mathfrak{f}}_{\max}]$ , where  $i < j$ ,  $\zeta^i - \zeta^j = \zeta^i \cdot (1 - \zeta^{j-i})$  is invertible due to Lemma 2. The expansion factor satisfying  $\gamma_S = 1$  is immediate. For the inverse-expansion factor, we have

$$\theta_S = \max_{i \neq j} \left\| \frac{1}{\zeta^i - \zeta^j} \right\|_{\infty} = \max_{i \neq j} \left\| \frac{1}{1 - \zeta_{i-j}} \right\|_{\infty} \leq \frac{\mathfrak{f}}{4\sqrt{2}}.$$

where the inequality is due to Lemma 1.  $\square$

### 4.3 Real Cyclotomics

We identify subtractive sets for real cyclotomic rings, i.e. the rings of integers of maximal real subfields of cyclotomic fields. The results over these rings mirror those for cyclotomic fields presented in Theorems 1 and 2.

**Theorem 3.** *Let  $\mathcal{R}^+ = \mathbb{Z}[\zeta_{\mathfrak{f}} + \zeta_{\mathfrak{f}}^{-1}]$  with  $\mathfrak{f} = p^k$ ,  $\mathfrak{f} > 4$ ,  $p$  prime. The set*

$$S := \{\mu_1^+, \dots, \mu_{(p+1)/2}^+\} \subseteq_{(p+1)/2} \mathcal{R}^+$$

*is subtractive, where  $\mu_i^+ = \mu_i + \bar{\mu}_i$  and  $\mu_i = (\zeta^i - 1)/(\zeta - 1)$  for  $i \in [(p+1)/2]$ , where  $\bar{\cdot}$  denotes the complex conjugate. Furthermore,  $\gamma_S \leq p$  and  $\theta_S \leq \frac{\mathfrak{f}^2}{8}$ .*

*Proof.* Observe that, for  $i > j$ ,

$$\mu_i - \mu_j = (1 + \zeta^{-j-i+1}) \cdot \frac{\zeta^i - \zeta^j}{\zeta - 1}.$$

The first factor is invertible if  $j + i - 1 \nmid \mathfrak{f}$ , which holds for distinct  $i, j \in [(p+1)/2]$ . The second factor is invertible due to Theorem 1. Hence,  $S$  is subtractive.

Since any  $c \in S$  is a sum of at most  $p$  roots of unity, we have  $\gamma_S \leq p$ . For  $\theta_S$ , we observe that

$$\frac{1}{\mu_i - \mu_j} = \frac{\zeta - 1}{(1 + \zeta^{-j-i+1}) \cdot (\zeta^i - \zeta^j)}.$$

Write  $\|\cdot\| = \|\sigma(\cdot)\|_{\infty}$ . By Lemma 1,

$$\theta_S \leq \left\| \frac{\zeta - 1}{(1 + \zeta^{-j-i+1}) \cdot (\zeta^i - \zeta^j)} \right\| \leq \|\zeta^i \cdot (\zeta - 1)\| \cdot \left\| \frac{1}{1 - \zeta^{j-i}} \right\| \cdot \left\| \frac{1}{1 + \zeta^{-j-i+1}} \right\| \leq \frac{\mathfrak{f}^2}{8}. \quad \square$$

**Theorem 4.** Let  $\mathcal{R}^+ = \mathbb{Z}[\zeta_f + \zeta_f^{-1}]$  with a non-prime-power conductor  $f$  having  $k \geq 2$  coprime prime-power factors  $(\hat{f}_i)_{i \in [k]}$ , i.e.  $f = \prod_{i \in [k]} \hat{f}_i$ . Write  $\hat{f}_{\max} := \max_{i \in [k]} \hat{f}_i$ . The set

$$S := \{\zeta^i + \zeta^{-i}\}_{\lfloor \frac{i}{\hat{f}_{\max}} \rfloor} \subseteq \lfloor \frac{i}{\hat{f}_{\max}} \rfloor \mathcal{R}^+,$$

is subtractive. Furthermore,  $\gamma_S \leq 2$  and  $\theta_S \leq \frac{f^2}{32}$ .

*Proof.* Consider  $c_i = \zeta^i + \zeta^{-i} \in S$  and  $c_j = \zeta^j + \zeta^{-j} \in S$  with  $i > j$ . Note that  $c_i - c_j = (\zeta^i + \zeta^{-i}) - (\zeta^j + \zeta^{-j}) = \zeta^{-i} \cdot (\zeta^{i+j} - 1) \cdot (\zeta^{i-j} - 1)$ . As,  $i + j, i - j \in \lfloor f/\hat{f}_{\max} \rfloor$ ,  $c_i - c_j$  is invertible in  $\mathcal{R}$  by Theorem 2.

The expansion factor satisfying  $\gamma_S \leq 2$  is immediate. Write  $\|\cdot\| = \|\sigma(\cdot)\|_\infty$ . For  $\theta_S$ , we observe that

$$\left\| \frac{1}{c_i - c_j} \right\| \leq \left\| \frac{1}{\zeta^{i+j} - 1} \right\| \cdot \left\| \frac{1}{\zeta^{i-j} - 1} \right\| \leq \left( \frac{f}{4\sqrt{2}} \right)^2 = \frac{f^2}{32},$$

where the inequality follows from Theorem 2.  $\square$

## 5 Succinct Arguments for Bounded-Norm Satisfiability

In this section, we assume that  $\mathcal{R}$  is either  $\mathcal{O}_K$  or  $\mathcal{O}_{K^+}$  which admit large enough subtractive sets, e.g. those constructed in Section 4. Let  $\mathcal{C}_{\mathcal{R}} \subset \mathcal{R}$  denote a fixed subtractive set with expansion factor  $\gamma$  and inverse-expansion factor  $\theta$ . Throughout, we mainly use the canonical 2-norm, and simply write  $\|\cdot\| := \|\sigma(\cdot)\|_2$ , unless specified. We use the shorthand notation  $\mathcal{R}_q^{n \times d^{\otimes \mu}} := ((\mathcal{R}_q^{1 \times d})^{\otimes \mu})^n$  for a matrix whose rows are elementary tensors. We also write  $\bar{\mathbf{Z}}$  (resp.  $\underline{\mathbf{Z}}$ ) to indicate the top (resp. bottom) half of a block matrix; the block dimension will be clear from the context. Lastly, we let  $\mathcal{C}_{\mathcal{R}_q} \subseteq \mathcal{R}_q^\times$  be obtained by taking a subfield of  $\mathcal{R}_q$  and removing 0. Note that  $\mathcal{C}_{\mathcal{R}}$  and  $\mathcal{C}_{\mathcal{R}_q}$  have the *invertible differences property* with respect to  $\mathcal{R}$  and  $\mathcal{R}_q$  respectively, i.e.  $\forall x \neq y \in \mathcal{C}_{\mathcal{R}}$  (resp.  $\mathcal{C}_{\mathcal{R}_q}$ ):  $x - y \in \mathcal{R}^\times$  (resp.  $\mathcal{R}_q^\times$ ).

We construct succinct arguments for proving that a short vector  $\mathbf{w}$  satisfy:

- $\mathcal{R}_q$ -linear elementary tensor relations, i.e.  $(\mathbf{g}_{\mu-1} \otimes \dots \otimes \mathbf{g}_0) \cdot \mathbf{w} = \mathbf{y} \bmod q$ ;
- a self-inner-product relation, i.e.  $t = \langle \mathbf{w}, \alpha(\mathbf{w}) \rangle_{\mathcal{R}} = \sum_{i=0}^{\mu-1} w_i \cdot \alpha(w)_i \in \mathcal{R}_q$ ; where  $\alpha \in \{\text{id}, \bar{\cdot}\}$  is either the identity map or the complex conjugate; and
- a norm bound  $\|\mathbf{w}\| \leq \beta$ .

In Table 1 on page 24, we provide an overview of parameters for correctness and relaxed knowledge soundness.

### 5.1 The (principal) relation $\Xi^{\text{lin}}$

We begin by defining the relation  $\Xi^{\text{lin}}$  and outline how protocols reduce instances in this relation to other instances. This relation serves as the principal building block for further protocols.

**Basic (single-block) relation.** We define our central relation(s) over the ring  $\mathcal{R}$ , modulo  $q$ , for witness dimension  $m = d^\mu$ . In fact, there are two central relations:  $\Xi^{\text{lin}}$  for correctness; and  $\Xi^{\text{lin}\vee\text{sis}}$  for relaxed knowledge soundness. We define both at once, so that  $\Xi^{\text{lin}\vee\text{sis}} \supseteq \Xi^{\text{lin}}$  contains all **highlighted** parts *additionally*. Let

$$\Xi_{\mathcal{R}, q, m, n^{\text{out}}, \mu, \beta, \beta^{\text{sis}}}^{\text{lin}\vee\text{sis}} := \left\{ \begin{array}{l} ((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}): \\ \mathbf{H} \in \mathcal{R}_q^{n^{\text{out}} \times n}; \mathbf{F} \in \mathcal{R}_q^{n \times d^{\otimes \mu}} \subseteq \mathcal{R}_q^{n \times m}; \mathbf{y} \in \mathcal{R}_q^{n^{\text{out}}} \\ \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta \\ \mathbf{H}\mathbf{F}\mathbf{w} = \mathbf{H}\mathbf{y} \bmod q \end{array} \right\} \quad \text{or} \quad \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta^{\text{sis}} \\ \bar{\mathbf{H}}\mathbf{F}\mathbf{w} = \mathbf{0}_{\bar{n}} \bmod q \end{array} \right\} \end{array} \right\}$$

where we *always assume* that  $\mathbf{H}$  has the block structure<sup>13</sup>

$$\mathbf{H} = \begin{pmatrix} \bar{\mathbf{H}} \\ \underline{\mathbf{H}} \end{pmatrix} \in \mathcal{R}_q^{n^{\text{out}} \times n} \quad \text{where} \quad \bar{\mathbf{H}} = (\mathbf{I}_{\bar{n}} \mathbf{0}) \in \mathcal{R}_q^{\bar{n} \times n} \quad \text{and} \quad \underline{\mathbf{H}} \in \mathcal{R}_q^{n^{\text{out}} \times n} \quad (2)$$

<sup>13</sup>This can be marginally relaxed: As long as there is an invertible  $\mathbf{X} \in \mathcal{R}^{n^{\text{out}} \times n^{\text{out}}}$  such that  $\mathbf{X}\mathbf{H}$  has this block structure, we can replace the claim  $(\mathbf{H}, \mathbf{F}, \mathbf{y})$  with the equivalent claim  $(\mathbf{X}\mathbf{H}, \mathbf{F}, \mathbf{X}\mathbf{y})$  which has the block structure our protocols require.

Similarly, write  $\bar{\mathbf{y}} \in \mathcal{R}_q^{\bar{n}}$  and  $\underline{\mathbf{y}} \in \mathcal{R}_q^{\underline{n}}$  for the  $\bar{n}$  top (resp.  $\underline{n}$  bottom) rows of  $\mathbf{y}$ .

*Remark 1 (Notational conventions).* We often omit irrelevant parameters in  $\Xi^{\text{lin}}$  and similar relations. Especially all fixed parameters in our protocols, which are  $\mathcal{R}$ ,  $q$ ,  $\bar{n}$ ,  $\beta^{\text{sis}}$ . For example, for parameterised relation like  $\Xi_{\mathcal{R},q,x,y}$ , we write  $\Xi_{x=f(\xi)}$  for  $\Xi_{\mathcal{R},q,f(\xi),z}$  or even just  $\Xi_{f(\xi)}$  if  $x = f(\xi)$  is clear from the context. Also, we fix  $d$  and always set  $m = d^\mu$ . As such, we often omit  $d$  and  $\mu$ .

Clearly, relation  $\Xi^{\text{lin}}$  asserts that the witness  $\mathbf{w}$  has norm  $\|\mathbf{w}\| \leq \beta$ . For the linear relation, let us first assume that  $\mathbf{H} = \mathbf{I}_n$  is an identity matrix. In this case, the relation asserts that  $\mathbf{F}\mathbf{w} = \mathbf{y}$  holds over  $\mathcal{R}_q$ . The matrix  $\mathbf{F}$  is structured, namely each row  $\mathbf{f}$  is an elementary tensor in  $\mathcal{R}_q^{1 \times d^{\otimes \mu}}$ , i.e.  $\mathbf{f} = \mathbf{g}_{\mu-1} \otimes \dots \otimes \mathbf{g}_0$  for  $\mathbf{g}_i = (g_{i,0}, \dots, g_{i,d-1}) \in \mathcal{R}_q^{1 \times d}$ .

For  $\Xi^{\text{lin}\text{v}\text{sis}}$ , we relax these assertions by introducing the **highlighted** OR-part, which captures a break of some underlying cryptographic assumption, e.g. a break of the vSIS assumption [CLM23] (Definition 1). For this,  $\mathbf{F} = \mathbf{H}\mathbf{F}$  will be the commitment key in a protocol. If the assumption is broken, then  $\Xi^{\text{lin}}$  may not be satisfied, hence the relaxed soundness relation  $\Xi^{\text{lin}\text{v}\text{sis}}$  is necessary.

Now, we further explain  $\mathbf{H}$ . The primary use of  $\mathbf{H}$  is to capture *random linear combination* of rows of  $\mathbf{F}$ . The block structure asserts that the top  $\bar{n}$  rows of  $\mathbf{F}$  are simply copied —  $\mathbf{F} = \mathbf{H}\mathbf{F}$  will correspond to the commitment key. Naively, our protocols would have communication costs linear in the number of rows of  $\mathbf{F}$ , but by using  $\mathbf{H}$ , we can compress this from  $n$  down to  $n^{\text{out}} = \bar{n} + \underline{n}$ . In prior works, one would simply (re)define  $\mathbf{F}$  as  $\mathbf{H}\mathbf{F}$  and  $\mathbf{y}$  as  $\mathbf{H}\mathbf{y}$ . However, to keep (*verifier-*)*succinctness*, we cannot do this: A (random) linear combination of elementary tensors is in general not an elementary tensor. However, our protocol crucially relies on the rows of  $\mathbf{F}$  being elementary tensor in order to apply FRI-style (verifier-succinct) folding of the statement. Therefore, we remember the (random) linear combinations of rows in  $\mathbf{H}$ , instead of carrying out the multiplication. Importantly, the *communication* of the protocol can indeed be compressed by applying  $\mathbf{H}$ . (Note there that the dimensions of  $\mathbf{H}$  and  $\mathbf{y}$  are in general much smaller than that of  $\mathbf{w}$ .)

**Reductions between  $\Xi^{\text{lin}}$ .** Our protocols reduce instances of  $\Xi_{m,\beta}^{\text{lin}}$  with different parameters, and we chain them to obtain our final split-and-fold protocol with intermediate norm checks. Primary protocols and parameters of interest are:

- (i)  $\Pi^{b\text{-decomp}}$ : Reduce one instance with bound  $\beta$  to many with bound  $b \ll \beta$ .
- (ii)  $\Pi^{\text{split}}$ : Reduce one instance with witness dimension  $m$  to many with  $m' = m/d$ .
- (iii)  $\Pi^{\text{fold}}$ : Reduce many instances with bounds  $\beta_i$  to one with  $\beta' = \gamma \sum_i \beta_i$ .
- (iv)  $\Pi^{\text{batch}}$ : Reduce one instance to another instance by randomly combining the last  $\underline{n}$  rows of  $\mathbf{H}$  and  $\mathbf{y}$  into a single one, so that  $n^{\text{out}} = \bar{n} + 1$ .

**Handling vSIS breaks.** Knowledge reductions can simply pass a  $\Xi^{\text{lin}\text{v}\text{sis}}$ -witness on as their extracted witness. Thus, we sometimes omit that discussion entirely.

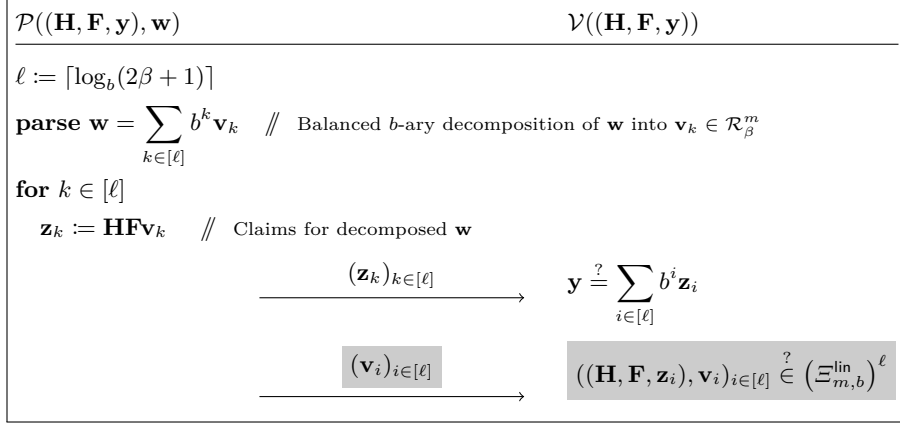
## 5.2 $\Pi^{b\text{-decomp}}$ : $b$ -ary Decomposition Knowledge Reduction

Let  $b \geq 1$  be an integer. The protocol  $\Pi^{b\text{-decomp}}$  (Fig. 4) is very simple: It takes a claim  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}) \in \Xi_{m,\beta}^{\text{lin}}$  and does a balanced  $b$ -ary decomposition of the witness  $\mathbf{w}$  with  $\|\mathbf{w}\| \leq \beta$  into  $\mathbf{w} = \sum_{i=0}^{\ell-1} b^i \mathbf{v}_i$ , where  $\mathbf{v}_i \in \mathcal{R}_b$  (hence  $\|\mathbf{v}_i\|_\infty \leq b/2$ ) and  $\ell = \lceil \log_b(2\beta + 1) \rceil$ . Then, appropriate claims  $\mathbf{z}_i = \mathbf{H}\mathbf{F}\mathbf{v}_i$  for the decomposed witness are computed, and the verifier makes sure the new claims imply the original one. Thus, the original statement is reduced to  $((\mathbf{H}, \mathbf{F}, \mathbf{z}_i), \mathbf{v}_i)_{i \in [\ell]}$ .

*Remark 2.* In protocol  $\Pi^{b\text{-decomp}}$ , we could apply the optimisation of not sending  $\mathbf{z}_0$ , and instead let the verifier compute the unique accepting  $\mathbf{z}_0$ , i.e. such that  $\mathbf{y} = \sum_{i \in [\ell]} b^i \mathbf{z}_i$ .

**Lemma 3.** *The protocol  $\Pi^{b\text{-decomp}}$  is a reduction from  $\Xi_\beta^{\text{lin}}$  to  $(\Xi_{\frac{1}{2}\sqrt{m}\varphi^{3/2}b}^{\text{lin}})^\ell$  for  $\ell = \lceil \log_b(2\beta + 1) \rceil$ . It is perfectly correct. It is perfectly relaxed knowledge sound from  $\Xi_{\frac{b^\ell-1}{b-1}\beta'}^{\text{lin}\text{v}\text{sis}}$  to  $(\Xi_{\beta'}^{\text{lin}\text{v}\text{sis}})^\ell$  if  $\frac{b^\ell-1}{b-1} \cdot \beta' \leq \beta^{\text{sis}}$ .*

*Note that  $\frac{b^\ell-1}{b-1} \cdot \beta' \leq 2b^{\ell-1}\beta' \leq 2\beta\beta'$ .*



**Fig. 4.** Protocol  $\Pi^{b\text{-decomp}}$ , a reduction of  $\Xi_{m,\beta}^{\text{lin}}$  to  $(\Xi_{m,b}^{\text{lin}})^\ell$  for  $\ell = \lceil \log_b(2\beta + 1) \rceil$ . As a *proof of knowledge*, Send the **marked** parts; as a *reduction not*,  $\Pi^{b\text{-decomp}}$  sends the **marked** parts only as a *proof* (but not *reduction*) of knowledge.

*Proof.* Perfect correctness of  $\Pi^{b\text{-decomp}}$  from  $\Xi_{m,\beta}^{\text{lin}}$  to  $(\Xi_{m,b}^{\text{lin}})^\ell$  is easy to see: By construction, each  $\mathbf{v}_i$  has  $\|\psi(\mathbf{v}_i)\|_\infty \leq b/2$ , and by Lemma 9 it follows that  $\|\mathbf{v}_i\| \leq \frac{1}{2}\sqrt{m}\varphi^{3/2}b$  as claimed. The linear equations  $\mathbf{H}\mathbf{F}\mathbf{v}_i = \mathbf{z}_i$  hold by definition for  $i \neq 0$  and by linearity for  $i = 0$ .

For relaxed knowledge soundness, observe that again by linearity, the original linear equation holds for  $\mathbf{w} = \sum_{i=0}^{\ell-1} b^i \mathbf{v}_i$ . For the norm, we have

$$\|\mathbf{w}\| \leq \sum_{i=0}^{\ell-1} b^i \|\mathbf{v}_i\| \leq \frac{b^\ell - 1}{b - 1} \cdot \beta'.$$

The OR-branch in  $\Xi^{\text{lin vsis}}$  is handled by letting the  $\mathbf{v}_i$  with  $\overline{\mathbf{H}}\mathbf{F}\mathbf{v}_i = \mathbf{0}$  be the extracted witness. If  $\frac{b^\ell - 1}{b - 1} \cdot \beta' \leq \beta^{\text{sis}}$ , this is a witness for  $\Xi^{\text{lin vsis}}$ .

For the note, we observe that by geometric series  $\frac{b^\ell - 1}{b - 1} \leq \frac{b^\ell}{b - (1 - 1/b)} = \frac{b}{b - 1} b^{\ell - 1} \leq 2b^{\ell - 1}$ . Moreover, by definition of  $\ell$ , we have  $b^{\ell - 1} < 2\beta + 1 \leq b^\ell$ . This yields the bounds.  $\square$

### 5.3 $\Pi^{\text{split}}$ : Witness Splitting Knowledge Reduction

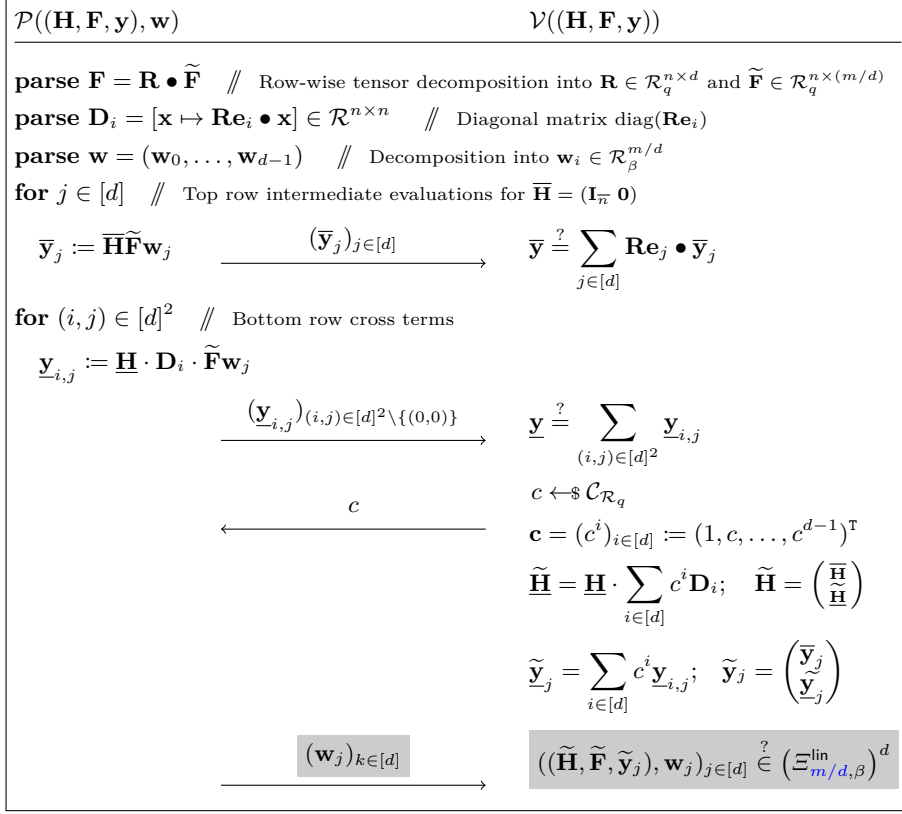
In Fig. 5 we describe protocol  $\Pi^{\text{split}}$  which takes a claim from  $\Xi_{m,\beta}^{\text{lin}}$  and splits it into  $d$  claims in  $\Xi_{m/d,\beta}^{\text{lin}}$ . That is, the number of claims in  $\Xi_{m,\beta}^{\text{lin}}$  (i.e. 1) grows by  $d$ -fold to  $d$ , but the witness dimension of each claim shrinks by  $d$ -fold to  $m/d$ . We explain the idea and correctness of the protocol below.

To split the witness, interpret  $\mathcal{R}^m$  as  $\mathcal{R}^{d \otimes \mu}$ , and split  $\mathbf{w} \in \mathcal{R}^m \cong \mathcal{R}^{d \otimes \mu}$  into  $\mathbf{w} = (\mathbf{w}_i)_{i \in [\mu]} = \sum_{i=0}^{\mu-1} \mathbf{e}_i \otimes \mathbf{w}_i$  where  $\mathbf{w}_i \in \mathcal{R}^{m/d} \cong \mathcal{R}^{d \otimes (\mu-1)}$  and  $\mathbf{e}_i \in \{0, 1\}^d$  is the  $i$ -th standard unit vector. Splitting  $\mathbf{w}$  like this is compatible with the row-wise tensor structure of  $\mathbf{F}$ . Let us take a closer look at this.

For simplicity, first consider a single row  $\mathbf{f} \in \mathcal{R}_q^{1 \times d \otimes \mu}$  of  $\mathbf{F}$ . By the elementary tensor structure of the row-vector  $\mathbf{f}$ , we can write it as  $\mathbf{f} = \mathbf{r} \otimes \tilde{\mathbf{f}} = (r_0 \cdot \tilde{\mathbf{f}}, \dots, r_{d-1} \cdot \tilde{\mathbf{f}}) = (\mathbf{f}_0, \dots, \mathbf{f}_{d-1})$  where  $\tilde{\mathbf{f}} \in \mathcal{R}_q^{1 \times d \otimes (\mu-1)}$ ,  $\mathbf{r} = (r_0, \dots, r_{d-1}) \in \mathcal{R}_q^{1 \times d}$ , and  $\mathbf{f}_i = r_i \cdot \tilde{\mathbf{f}}$ . Therefore,  $\mathbf{f} \cdot \mathbf{w} = \sum_{i \in [d]} \mathbf{f}_i \mathbf{w}_i = \sum_{i \in [d]} (\tilde{\mathbf{f}} \cdot \mathbf{w}_i) \cdot (\mathbf{r} \cdot \mathbf{e}_i^T) = \sum_{i \in [d]} r_i \tilde{\mathbf{f}} \cdot \mathbf{w}_i$ .

Now, consider any matrix  $\mathbf{F}$  with row-wise tensor structure and  $n$  rows, as in  $\Xi^{\text{lin}}$ . That is,  $\mathbf{F} \in \mathcal{R}_q^{n \times d \otimes \mu}$ . Observe that

$$\mathbf{F} = \begin{pmatrix} \mathbf{F}_{0,\bullet} \\ \vdots \\ \mathbf{F}_{n-1,\bullet} \end{pmatrix} = \begin{pmatrix} \mathbf{F}_{0,0} & \dots & \mathbf{F}_{0,d-1} \\ \vdots & & \vdots \\ \mathbf{F}_{n-1,0} & \dots & \mathbf{F}_{n-1,d-1} \end{pmatrix} = \begin{pmatrix} \mathbf{r}_0 \otimes \tilde{\mathbf{f}}_0 \\ \vdots \\ \mathbf{r}_{n-1} \otimes \tilde{\mathbf{f}}_{n-1} \end{pmatrix} \quad (3)$$



**Fig. 5.** Protocol  $\Pi^{\text{split}}$ , a reduction from  $\Xi_{m, \beta}^{\text{lin}}$  to  $(\Xi_{m/d, \beta}^{\text{lin}})^d$ .  $\Pi^{\text{split}}$  sends the marked parts only as a *proof* (but not *reduction*) of knowledge.

where  $\mathbf{F}_{i, \bullet}$  denotes the  $i$ -th row of  $\mathbf{F}$ , and  $\mathbf{F}_{i,j} \in \mathcal{R}_q^{1 \times d^{\otimes \mu}}$  the block of rows (the analogue of  $(\mathbf{f}_0, \dots, \mathbf{f}_{d-1})$  of the single-row case), and  $\tilde{\mathbf{f}}_i \in \mathcal{R}_q^{1 \times d^{\otimes (\mu-1)}}$  and  $\mathbf{r}_i \in \mathcal{R}_q^{1 \times d}$  are the analogues of  $\mathbf{r}$  and  $\tilde{\mathbf{f}}$  of the single-row case respectively. To ease notation, we define  $\mathbf{R} = (\mathbf{r}_i^\top)_{i \in [n]} \in \mathcal{R}_q^{n \times d}$  and  $\tilde{\mathbf{F}} = (\tilde{\mathbf{f}}_i)_{i \in [n]} \in \mathcal{R}_q^{n \times d^{\otimes (\mu-1)}}$ , and we write  $\mathbf{F} = \mathbf{R} \bullet \tilde{\mathbf{F}}$  for the row-wise tensor product<sup>14</sup> of  $\mathbf{R}$  and  $\tilde{\mathbf{F}}$  as seen in Eq. (3). In this notation,

$$\mathbf{F} \cdot \mathbf{w} = (\mathbf{R} \bullet \tilde{\mathbf{F}}) \cdot \left( \sum_{i=0}^{\mu-1} \mathbf{e}_i \otimes \mathbf{w}_i \right) = \sum_i \underbrace{(\mathbf{R}\mathbf{e}_i)}_{\in \mathcal{R}_q^n} \odot \underbrace{(\tilde{\mathbf{F}}\mathbf{w}_i)}_{=\mathbf{y}_i \in \mathcal{R}_q^n} \quad (4)$$

where we use the Hadamard product to multiply the vector  $\mathbf{R}\mathbf{e}_i$  with  $\mathbf{y}_i$  componentwise. Moreover, with  $\mathbf{D}_i := \text{diag}(\mathbf{r}_i)$ , we can rewrite (4) as

$$\mathbf{F} \cdot \mathbf{w} = \sum_i \mathbf{D}_i \cdot \tilde{\mathbf{F}}\mathbf{w}_i \quad (5)$$

With the above, we have derived a splitting protocol for the special case where  $\mathbf{H} = \mathbf{I}_n$  is the identity matrix: Simply send  $\tilde{\mathbf{y}}_i = \tilde{\mathbf{F}}\mathbf{w}_i$  for new statements  $(\mathbf{H}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}_i)$  and witnesses  $(\mathbf{w}_i)_{i \in [d]}$ .

When  $\mathbf{H}$  is not necessarily the identity, we must also handle the bottom part  $\underline{\mathbf{H}}$  of  $\mathbf{H}$ . To do so, our protocol (cf. Fig. 5) additionally sends cross terms, namely  $\mathbf{D}_i \cdot \tilde{\mathbf{F}}\mathbf{w}_j$  for  $i, j \in [d]$ , which are then randomly recombined.

*Remark 3.* In protocol  $\Pi^{\text{split}}$ , we could apply the optimisation of not sending  $\bar{\mathbf{y}}_j$  (assuming that  $\mathbf{R}\mathbf{e}_j \neq \mathbf{0}$ ) and  $\underline{\mathbf{y}}_{0,0}$ , and instead let the verifier compute the unique accepting  $\bar{\mathbf{y}}_j$  and  $\underline{\mathbf{y}}_{0,0}$ , i.e. such that  $\bar{\mathbf{y}} = \sum_{j \in [d]} \mathbf{R}\mathbf{e}_j \bullet \bar{\mathbf{y}}_j$  and  $\underline{\mathbf{y}} = \sum_{(i,j) \in [d]^2} \underline{\mathbf{y}}_{i,j}$ .

<sup>14</sup>(and more general forms, as block Kronecker product and Khatri–Rao product).



**Lemma 4.** *The protocol  $\Pi^{\text{split}}$  is a reduction of knowledge from  $\Xi_{m,\mu,\beta}^{\text{lin}}$  to  $(\Xi_{m/d,\mu-1,\beta}^{\text{lin}})^d$ . It is perfectly correct. It is relaxed knowledge sound from  $\Xi_{m,\mu,\sqrt{d}\beta}^{\text{lin}\text{vis}}$  to  $(\Xi_{m/d,\mu-1,\beta}^{\text{lin}})^d$  with  $d$ -special sound extraction and knowledge error  $\kappa = (d-1)/|\mathcal{C}_{\mathcal{R}_q}|$  if  $2\beta \leq \beta^{\text{sis}}$ .*

*Proof.* Perfect correctness of  $\Pi^{\text{split}}$  from  $\Xi_{m,\mu,\beta}^{\text{lin}}$  to  $\Xi_{m/d,\mu-1,\beta}^{\text{lin}}$  is straightforward for the top rows: Since  $\bar{\mathbf{H}} = (\mathbf{I}_{\bar{n}} \mathbf{0})$ , we have  $\bar{\mathbf{F}} = \bar{\mathbf{H}}\mathbf{F}$  are just the  $\bar{n}$  top rows of  $\mathbf{F}$ , and similar for  $\tilde{\mathbf{F}}$ , and thus by our discussion before and some renaming (using  $\bar{\mathbf{F}}$  instead of  $\mathbf{F}$  makes  $\bar{\mathbf{H}}$  the identity) we know that the top part is perfectly correct. For the bottom rows, correctness follows essentially from Eqs. (4) and (5) which asserts that

$$\tilde{\mathbf{H}}\tilde{\mathbf{F}}\mathbf{w}_j = \sum_{i,j \in [d]} c_i \underline{\mathbf{H}} \cdot \mathbf{D}_i \cdot \tilde{\mathbf{F}}\mathbf{w}_j = \sum_{i,j \in [d]} c_i \underline{\mathbf{y}}_{i,j} = \tilde{\mathbf{y}}_j.$$

For relaxed knowledge soundness, we argue through  $d$ -special soundness. So we have  $d$  related accepting transcripts for challenge vectors  $\mathbf{c}^{(k)}$  with witnesses  $\mathbf{w}_j^{(k)}$  which satisfy  $((\tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}_j), \mathbf{w}_j) \in \Xi_{m/d,\mu-1}^{\text{lin}\text{vis}}$ .

*Step 1 (top rows):* Let us first consider the top rows (and any single transcript): Here, it is straightforward to see that

$$\mathbf{w}^{(k)} = \begin{pmatrix} \mathbf{w}_0^{(k)} \\ \vdots \\ \mathbf{w}_{d-1}^{(k)} \end{pmatrix} \text{ satisfies } \bar{\mathbf{H}}\mathbf{F}\mathbf{w}^{(k)} = \bar{\mathbf{H}} \sum_i \mathbf{D}_i \tilde{\mathbf{F}}\mathbf{w}_i^{(k)} = \sum_i \bar{\mathbf{y}}_i = \bar{\mathbf{y}}$$

for all  $k \in [d]$  by construction (and using  $\bar{\mathbf{H}} = (\mathbf{I}_{\bar{n}} \mathbf{0})$ ). Thus, we trivially and unconditionally find a witness for the top rows. Clearly,  $\|\mathbf{w}\|_2 \leq \sqrt{d}\beta'$ .

Moreover, by looking at the top rows, we see that: Either, there is a unique  $\mathbf{w}_j$  over all transcripts, i.e.  $\mathbf{w}_j^{(k)} = \mathbf{w}_j^{(k')}$  for all  $k, k' \in [d]$ . Or, there is a *non-zero* difference  $\mathbf{v}_j = \mathbf{w}_j^{(k)} - \mathbf{w}_j^{(k')}$  of norm at most  $2\beta$ , such that  $\tilde{\mathbf{H}}\tilde{\mathbf{F}}\mathbf{v}_j = \mathbf{0}$ , and thus  $\bar{\mathbf{H}}\mathbf{F}(\mathbf{v}_j, \mathbf{0}, \dots, \mathbf{0}) = \mathbf{0}$  is a witness for the OR-branch in  $\Xi^{\text{lin}\text{vis}}$  (of norm at most  $2\beta$ ). Hence, from now on, we assume all transcripts contain the same  $\mathbf{w}_j = \mathbf{w}_j^{(k)}$  for all  $k \in [d]$ .

*Step 2 (bottom rows):* Now, we consider the bottom row, with an arbitrary  $\underline{\mathbf{H}}$ . Towards showing  $d$ -special soundness, define for  $i \in [0, \mu-1]$  and  $j \in [0, d-1]$  the shorthand

$$\mathbf{z}_{i,j} = \mathbf{D}_i \tilde{\mathbf{F}}\mathbf{w}_j.$$

As the first step, we show that  $\underline{\mathbf{H}}\mathbf{z}_{i,j} = \underline{\mathbf{y}}_{i,j}$  for all  $i, j$ . Towards this, we rewrite the verifier's checks as

$$\underline{\mathbf{H}} \cdot \mathbf{Z}_j \cdot \mathbf{c}^{(k)} = \sum_i \underline{\mathbf{H}} \cdot (\mathbf{z}_{i,j})_i c_i^{(k)} = \sum_i (\underline{\mathbf{y}}_{i,j})_i c_i^{(k)} = \underline{\mathbf{Y}}_j \cdot \mathbf{c}^{(k)}$$

where  $\mathbf{Z}_j = (\mathbf{z}_{0,j}, \dots, \mathbf{z}_{d-1,j})$  and likewise for  $\underline{\mathbf{Y}}_j$ . From  $d$  distinct challenges, we assemble a (Vandermonde) matrix  $\mathbf{C} = (\mathbf{c}^{(0)}, \dots, \mathbf{c}^{(d-1)})$ . Since  $\mathcal{C}_{\mathcal{R}_q}$  has the invertibility of differences property,  $\mathbf{C}$  is invertible over  $\mathcal{R}_q$ , and therefore

$$\underline{\mathbf{H}} \cdot \mathbf{Z}_j \cdot \mathbf{C} = \underline{\mathbf{Y}}_j \cdot \mathbf{C} \implies \underline{\mathbf{H}} \cdot \mathbf{Z}_j = \underline{\mathbf{Y}}_j.$$

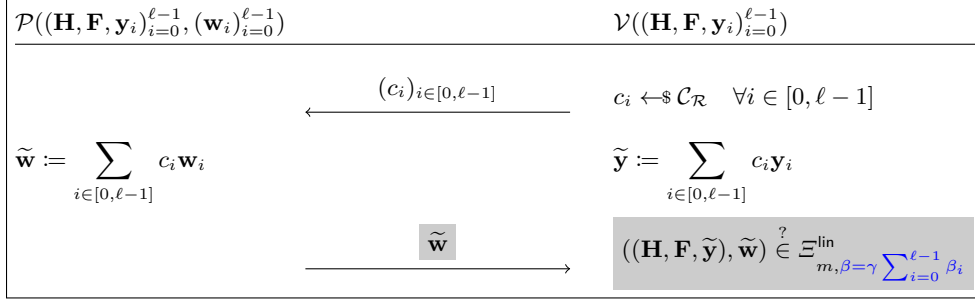
Thus  $\underline{\mathbf{H}}\mathbf{z}_{i,j} = \underline{\mathbf{y}}_{i,j}$  for all  $i, j$  as claimed. Then we see that

$$\underline{\mathbf{y}} = \sum_i \underline{\mathbf{y}}_{i,i} = \underline{\mathbf{H}} \sum_i \mathbf{z}_{i,i} = \underline{\mathbf{H}} \sum_i \mathbf{D}_i \tilde{\mathbf{F}}\mathbf{w}_i = \underline{\mathbf{H}}\mathbf{F}\mathbf{w}$$

and therefore,  $\mathbf{w}$  (as assembled in Step 1) is a witness for the bottom rows as well.

*Step 3 (OR-branch):* Finally consider the OR-branch in  $\Xi^{\text{lin}\text{vis}}$ . If  $\tilde{\mathbf{H}}\tilde{\mathbf{F}}\mathbf{v}_j = \mathbf{0}$ , we simply let  $\mathbf{v}_j$  be the extracted witness and note that  $\|\mathbf{v}_j\| \leq \beta \leq \beta^{\text{sis}}$ .  $\square$

*Remark 4.* Protocol  $\Pi^{\text{split}}$  can be optimised. For example, suppose that  $\mathbf{r}_0 = \mathbf{R}\mathbf{e}_0$  has no zero component. Then instead of sending  $\bar{\mathbf{y}}_0$ , we can compute it as  $\mathbf{D}_0^{-1} \sum_{i \in [d] \setminus \{0\}} \mathbf{D}_i \bar{\mathbf{y}}_i$  because no other choice satisfies the verifier's check. Similarly, we can omit  $\underline{\mathbf{y}}_{0,0}$ . For arbitrary  $\mathbf{R}$ , in each row there is some  $i$  such that  $\mathbf{R}\mathbf{e}_i$  is not zero (else  $\mathbf{F}$  has a zero row, which is useless), so a more complex variant of this optimization always applies, saving  $n^{\text{out}}$   $\mathcal{R}_q$ -elements of communication. Moreover, whenever  $\underline{\mathbf{H}}$  has structured rows (e.g., contains (permuted) identity submatrices, etc.), application specific optimisations may apply.



**Fig. 6.** Protocol  $\Pi^{\text{fold}}$  folds multiple instances of  $\Xi^{\text{lin}}$  with the same  $(\mathbf{H}, \mathbf{F})$  into one.  $\Pi^{\text{fold}}$  sends the marked parts only as a *proof* (but not *reduction*) of knowledge.

#### 5.4 $\Pi^{\text{fold}}$ : Fold Knowledge Reduction

In Fig. 6, we present the protocol  $\Pi^{\text{fold}}$ , which is a simple batch verification for many statements of the same type. It takes  $\ell$  instances of  $((\mathbf{H}, \mathbf{F}, \mathbf{y}_i), \mathbf{w}_i)_{i \in [\ell]}$  of  $\Xi_{m, \beta}^{\text{lin}}$  with the same  $(\mathbf{H}, \mathbf{F})$ , and produces a random linear combination  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \tilde{\mathbf{w}})$  as output, with increased norm bounds.

**Lemma 5.** *The protocol  $\Pi^{\text{fold}}$  is a reduction of knowledge from  $(\Xi_{\beta}^{\text{lin}})^{\ell}$  to  $\Xi_{\gamma\ell\beta}^{\text{lin}}$ . It is perfectly correct. It is  $\ell$ -CWSS and relaxed knowledge sound from  $(\Xi_{2\theta\beta'}^{\text{lin}\vee\text{sis}})^{\ell}$  to  $\Xi_{\beta'}^{\text{lin}\vee\text{sis}}$  if  $\beta' \leq \beta^{\text{sis}}$ .*

*Proof.* For perfect correctness, it is clear that  $\tilde{\mathbf{w}}$  satisfies  $\mathbf{H}\mathbf{F}\tilde{\mathbf{w}} = \tilde{\mathbf{y}}$  by construction. Moreover,  $\|\tilde{\mathbf{w}}\| \leq \sum_{i=1}^{\ell-1} \|c_i \mathbf{w}_i\| \leq \sum_{i=1}^{\ell-1} \gamma \|\mathbf{w}_i\| \leq \gamma\ell\beta$ , thus the norm is also within bounds and correctness follows.

For relaxed knowledge soundness, through  $\ell$ -CWSS we are given  $\ell+1$  accepting transcripts  $\tilde{\mathbf{w}}_0, \dots, \tilde{\mathbf{w}}_{\ell}$ , for challenges  $(c_0^{(i)}, \dots, c_{\ell-1}^{(i)})$  where  $\mathbf{c}^{(i)}$  and  $\mathbf{c}^{(\ell)}$  differ exactly in coordinate  $i \in \{0, \dots, \ell-1\}$ . We can now subtract the accepting equations to obtain

$$\mathbf{H}\mathbf{F}(\tilde{\mathbf{w}}_i - \tilde{\mathbf{w}}_{\ell}) = (\mathbf{c}^{(i)} - \mathbf{c}^{(\ell)})\mathbf{y} = (c_i^{(i)} - c_i^{(\ell)})\mathbf{y}_i$$

and thus, setting

$$\mathbf{w}_i := \frac{1}{c_i^{(i)} - c_i^{(\ell)}} (\tilde{\mathbf{w}}_i - \tilde{\mathbf{w}}_{\ell})$$

is a witness in  $\mathcal{R}_q^m$  satisfying  $\mathbf{H}\mathbf{F}\mathbf{w}_i = \mathbf{y}_i$ ; Here, we use the subtractive set property of  $\mathcal{C}_{\mathcal{R}}$ . Moreover, we have

$$\|\mathbf{w}_i\| = \left\| \frac{1}{c_i^{(i)} - c_i^{(\ell)}} (\tilde{\mathbf{w}}_i - \tilde{\mathbf{w}}_{\ell}) \right\| \leq \theta \cdot 2 \cdot \beta'$$

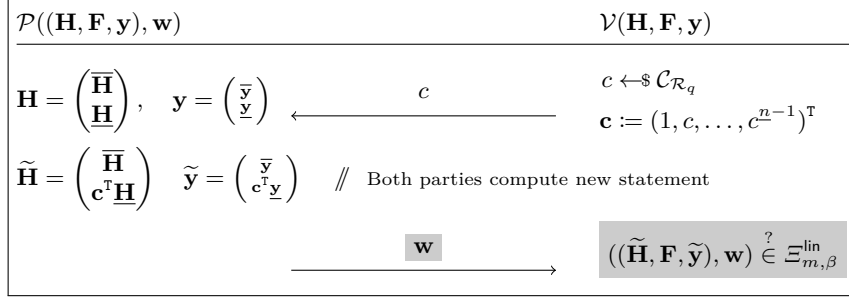
by definition of the inverse-expansion factor for  $\mathcal{C}_{\mathcal{R}}$  and the bounds on all  $\tilde{\mathbf{w}}_j$ .

Finally, the OR-branch in  $\Xi^{\text{lin}\vee\text{sis}}$  is handled by letting  $\mathbf{w}_i$  equal  $\tilde{\mathbf{w}}$  for all  $i$ ; obviously,  $\overline{\mathbf{H}}\mathbf{F}\mathbf{w}_i = \mathbf{0}$  holds and  $\|\mathbf{w}_i\| \leq \beta' \leq \beta^{\text{sis}}$ . This completes the proof.  $\square$

#### 5.5 $\Pi^{\text{batch}}$ : Batch-Rows Knowledge Reduction

The protocol  $\Pi^{\text{batch}}$  (Fig. 7) is a protocol to batch the claims along multiple rows into fewer rows of claims. This is done by a random linear combination of the rows in question. This protocol maps an instance  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w})$  of  $\Xi_{m, \beta}^{\text{lin}}$  to an instance  $((\overline{\mathbf{H}}, \mathbf{F}, \tilde{\mathbf{y}}), \mathbf{w})$ , where the dimension of  $\tilde{\mathbf{y}}$  is smaller. We describe it in more detail: Let  $n^{\text{out}} = \bar{n} + \underline{n}$ . Then  $\Pi^{\text{batch}}$  keeps the top  $\bar{n}$  rows  $\bar{\mathbf{y}}$  of  $\mathbf{y}$  (resp.  $\overline{\mathbf{H}}$  of  $\mathbf{H}$ , and thus of  $\mathbf{H}\mathbf{F}$ ) unchanged. But the bottom  $\underline{n}$  rows are linearly combined into a single row. For this,  $\mathbf{H}$  and  $\mathbf{y}$  are split into top and bottom half, and the bottom half is multiplied by a vector  $\mathbf{c}$  consisting of powers of  $c \leftarrow \mathcal{C}_{\mathcal{R}_q}$ . Both parties then update the statement suitably.

**Lemma 6.** *The protocol  $\Pi^{\text{batch}}$  is a reduction of knowledge from  $\Xi_{n^{\text{out}}, \beta}^{\text{lin}}$  to  $\Xi_{\bar{n}+1, \beta}^{\text{lin}}$ . It is perfectly correct. It is relaxed knowledge sound from  $\Xi_{n^{\text{out}}, \beta'}^{\text{lin}\vee\text{sis}}$  to  $\Xi_{\bar{n}+1, \beta'}^{\text{lin}\vee\text{sis}}$  with knowledge error  $\kappa = \frac{\bar{n}-1}{|\mathcal{C}_{\mathcal{R}_q}|} \leq \frac{\bar{n}}{|\mathcal{C}_{\mathcal{R}_q}|}$  if  $2\beta' \leq \beta^{\text{sis}}$ . Extraction requires two uniformly distributed transcripts (Section 3.2).*



**Fig. 7.** Protocol  $\Pi^{\text{batch}}$  reduces an instance of  $\Xi^{\text{lin}}$  to another in instance with fewer rows by batching to last  $\underline{n}$  of  $\mathbf{H}$ .  $\Pi^{\text{batch}}$  sends the marked parts only as a *proof* (but not *reduction*) of knowledge.

*Proof.* The correctness of this protocol is straightforward by linearity. For (knowledge) soundness, we rely on the Schwartz–Zippel lemma over  $\mathcal{C}_{\mathcal{R}_q}$ , which we recall is almost a subfield  $\mathbb{F}$  of  $\mathcal{R}_q$  except that 0 is missing. The lemma states that, for any degree- $d$  non-zero polynomial over  $\mathbb{F}$ , the probability that the polynomial evaluates to zero at a uniformly random point chosen from  $\mathcal{C}_{\mathcal{R}_q}$  is at most  $d \cdot |\mathcal{C}_{\mathcal{R}_q}|^{-1}$ . To translate this upper bound into a *knowledge* error, observe the following: If  $\mathcal{A}$  succeeds for 2 challenges, then the first transcript fixes some  $\mathbf{w}_1$  which satisfies  $((\tilde{\mathbf{H}}_1, \mathbf{F}, \tilde{\mathbf{y}}_1), \mathbf{w}_1) \in \Xi^{\text{lin}}$ . Suppose  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}_1) \notin \Xi^{\text{lin}}$ , i.e.  $\mathbf{w}_1$  is not a witness for the original statement. Then we observe that *at most* a fraction of  $\kappa = \frac{n-1}{|\mathcal{C}_{\mathcal{R}_q}|}$  challenges can be accepting for  $\mathbf{w}_1$  (by Schwartz–Zippel). In other words, if  $\mathcal{A}$  succeeds with probability  $\epsilon$ , then with probability at least  $\epsilon - \kappa$  the 2-transcript extractor successfully outputs two transcripts where the responses  $\mathbf{w}_1$  and  $\mathbf{w}_2$  differ.<sup>15</sup> Now,  $\mathbf{v} = \mathbf{w}_1 - \mathbf{w}_2$  is a non-zero preimage  $\overline{\mathbf{H}}\mathbf{F}\mathbf{v} = \mathbf{0}$  of norm at most  $2\beta' \leq \beta^{\text{sis}}$ , i.e. the OR-branch of  $\Xi^{\text{lin}\vee\text{sis}}$ .  $\square$

## 5.6 $\Pi^{\text{split}\&\text{fold}}$ : Split-and-Fold

We describe protocol  $\Pi^{\text{split}\&\text{fold}}$  which reduces the size of the witness. While it is perfectly correct, the extracted relation suffers a (small) growth in norm, i.e. it is only relaxed knowledge sound. The protocol  $\Pi^{\text{split}\&\text{fold}}$  proceeds as follows:

- (1)  $\Xi_{m, \beta_0}^{\text{lin}} \xrightarrow{\Pi^{b\text{-decomp}}} (\Xi_{m, b}^{\text{lin}})^\ell$  (Reduce norms by  $b$ -ary decomposition.)
- (2)  $(\Xi_{m, b}^{\text{lin}})^\ell \xrightarrow{\Pi^{\text{split}}} (\Xi_{m/d, b}^{\text{lin}})^{\ell d}$  (Split statements into smaller ones.)
- (3)  $(\Xi_{m/d, b}^{\text{lin}})^{\ell d} \xrightarrow{\Pi^{\text{fold}}} \Xi_{m/d, \beta_3}^{\text{lin}}$  (Fold statements into single one.)

The goal behind  $\Pi^{\text{split}\&\text{fold}}$  is to reduce the statement size by applying  $\Pi^{\text{split}}$  and then  $\Pi^{\text{fold}}$ . However, doing just this increases the norm the folded witness. To avoid this,  $\Pi^{\text{split}\&\text{fold}}$  first applies a  $b$ -ary decomposition which decreases the norm of  $\mathbf{w}$  sufficiently, with proper parameters, we get  $\beta_3 = \beta_0$  again. Correctness of this protocol is straightforward to show. Relaxed knowledge soundness also follows from relaxed soundness of the building blocks. Unfortunately, the norm of the extracted witness *does* grow, no matter how the parameters are chosen.

*Remark 5.* When applying  $\Pi^{\text{split}}$  to the product relation  $(\Xi_{m, b}^{\text{lin}})^\ell$ , it is crucial that a single challenge is reused among all instances. This ensures that if all  $(\mathbf{H}_i, \mathbf{F}_i)$  were identical before splitting, they still are identical after splitting.

**Theorem 5.** *Let  $\beta_0, b \in \mathbb{N}$  (suitably chosen). Then  $\Pi^{\text{split}\&\text{fold}}$  is perfectly correct from  $\Xi_{m, \beta_0}^{\text{lin}}$  to  $\Xi_{m/d, \beta_3}^{\text{lin}}$ , where  $\beta_3 = \gamma \ell d \frac{1}{2} \sqrt{m} \varphi^{3/2} \beta_0$  with  $\ell = \lceil \log_b(2\beta_0 + 1) \rceil$ . It is relaxed knowledge sound from  $\Xi_{m, \beta'_0}^{\text{lin}\vee\text{sis}}$  to  $\Xi_{m/d, \beta'_3}^{\text{lin}\vee\text{sis}}$  if  $\beta'_0 = 4\sqrt{d}\theta\beta_0\beta'_3 \leq \beta^{\text{sis}}$ . The knowledge error is at most  $\frac{\ell d}{|\mathcal{C}_{\mathcal{R}}|} + \frac{d-1}{|\mathcal{C}_{\mathcal{R}_q}|}$ .*

*Proof.* For correctness, the change in dimensions are clear. We write  $\beta_i$  for the bound after the  $i$ -th step, starting out with  $\beta_0$  and finishing with  $\beta_3$ . Let  $\ell = \lceil \log_b(2\beta_0 + 1) \rceil$  as in Lemma 3. We have

<sup>15</sup>We exploit that the challenges are *uniformly* distributed (conditioned on accepting).

- $\beta_1 = \frac{1}{2}\sqrt{m}\varphi^{3/2}b$  by Lemma 3.
- $\beta_2 = \beta_1$  by Lemma 4.
- $\beta_3 = \gamma\ell d\beta_2$  by Lemma 5.

Overall, we find

$$\beta_3 = \gamma\ell d\frac{1}{2}\sqrt{m}\varphi^{3/2}b.$$

For knowledge soundness, we work in the opposite direction. Starting out from a witness in  $\Xi_{m/d,\beta'_3}^{\text{lin}\vee\text{sis}}$ , we follow the bounds back to the start. We have

- $\beta'_2 = 2\theta\beta'_3$  by Lemma 5.
- $\beta'_1 = \sqrt{d}\beta'_2$  by Lemma 4.
- $\beta'_0 = 2\beta_0\beta'_1$  by Lemma 3.

where  $\beta_0$  is the fixed choice for  $\Pi^{b\text{-decomp}}$  (w.r.t. correctness). This yields overall

$$\beta'_0 = 4\sqrt{d}\theta\beta_0\beta'_3.$$

The knowledge error is  $\frac{d-1}{|\mathcal{C}_{\mathcal{R}_q}|}$  for  $\Pi^{b\text{-decomp}}$  and  $\Pi^{\text{split}}$ , and it is  $\frac{\ell d}{|\mathcal{C}_{\mathcal{R}}|}$  for  $\Pi^{\text{fold}}$ . (For completeness, we note that since  $\beta'_0$  is the largest bound, all intermediate bounds satisfy  $\beta'_i \leq \beta^{\text{sis}}$  as well, so the extractors work.)  $\square$

## 5.7 $\Pi^{\text{norm}}$ , $\Pi^{\text{norm}+}$ , $\Pi^{\text{ip}}$ , $\Pi^{\text{ip}+}$ : Norm and Inner Product Checks

To restrain the norm growth of the extracted witness, we introduce norm checks. First we present the “core” norm check protocol  $\Pi^{\text{norm}}$ , which handles the interesting part of the norm check by reducing the norm relation  $\Xi^{\text{norm}}$ , to multiple  $\Xi^{\text{lin}}$  relations. We then compose  $\Pi^{\text{norm}}$  with  $\Pi^{\text{batch}}$  and  $\Pi^{\text{fold}}$  to yield the “full” norm check protocol  $\Pi^{\text{norm}+}$ , which reduces the norm relation to a single  $\Xi^{\text{lin}}$  relation. At the core of  $\Pi^{\text{norm}}$  is a mechanism for checking the trace of an inner product. By removing the trace operation, we obtain similar protocols  $\Pi^{\text{ip}}$  and  $\Pi^{\text{ip}+}$  for proving inner product relations  $\Xi^{\text{ip}}$ . The relations  $\Xi^{\text{norm}}$  and  $\Xi^{\text{ip}}$ , as well as their variants  $\Xi^{\text{norm}\vee\text{sis}}$  and  $\Xi^{\text{ip}\vee\text{sis}}$ , are defined as follows.

$$\Xi_{\mathcal{R},q,m,n^{\text{out}}}^{\text{norm}\vee\text{sis}} := \left\{ \begin{array}{l} ((\mathbf{H}, \mathbf{F}, \mathbf{y}, \nu), \mathbf{w}) : \\ \mathbf{H} \in \mathcal{R}_q^{n^{\text{out}} \times n}, \mathbf{F} \in \mathcal{R}_q^{n \times d^{\otimes \mu}} \subseteq \mathcal{R}_q^{n \times m}; \mathbf{y} \in \mathcal{R}_q^{n^{\text{out}}}, \nu \leq \beta \\ \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \nu \\ \mathbf{H}\mathbf{F}\mathbf{w} = \mathbf{H}\mathbf{y} \text{ mod } q \end{array} \right\} \text{ or } \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta^{\text{sis}} \\ \overline{\mathbf{H}\mathbf{F}\mathbf{w}} = \mathbf{0}_{\bar{n}} \text{ mod } q \end{array} \right\} \end{array} \right\},$$

$$\Xi_{\mathcal{R},q,m,n^{\text{out}}}^{\text{ip}\vee\text{sis}} := \left\{ \begin{array}{l} ((\mathbf{H}, \mathbf{F}, \mathbf{y}, t), \mathbf{w}) : \\ \mathbf{H} \in \mathcal{R}_q^{n^{\text{out}} \times n}, \mathbf{F} \in \mathcal{R}_q^{n \times d^{\otimes \mu}} \subseteq \mathcal{R}_q^{n \times m}; \mathbf{y} \in \mathcal{R}_q^{n^{\text{out}}}, t \in \mathcal{R}_q \\ \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta \\ \mathbf{H}\mathbf{F}\mathbf{w} = \mathbf{H}\mathbf{y} \text{ mod } q \\ \langle \mathbf{w}, \alpha(\mathbf{w}) \rangle_{\mathcal{R}} = t \text{ mod } q \end{array} \right\} \text{ or } \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta^{\text{sis}} \\ \overline{\mathbf{H}\mathbf{F}\mathbf{w}} = \mathbf{0}_{\bar{n}} \text{ mod } q \end{array} \right\} \end{array} \right\}.$$

Note that, compared to  $\Xi^{\text{lin}}$ , the norm relation  $\Xi^{\text{norm}}$  differs in that the witness norm bound  $\nu \leq \beta$  is given as part of the statement, and a stricter norm relation  $\|\mathbf{w}\| \leq \nu$  is checked. Similarly,  $\Xi^{\text{ip}}$  differs from  $\Xi^{\text{lin}}$  in that the statement additionally includes an inner product value  $t$ , and the witness additionally satisfies an inner product relation. Furthermore, we note that  $\Xi^{\text{ip}}$  is parametrised by  $\alpha \in \{\text{id}, \overline{\text{id}}\}$  which is either the identity or complex conjugate, controlling which type of inner product is being considered.

**The core protocols  $\Pi^{\text{norm}}$  and  $\Pi^{\text{ip}}$ .** The protocols  $\Pi^{\text{norm}}$ ,  $\Pi^{\text{ip}}$  for  $\alpha = \text{id}$  and  $\Pi^{\text{ip}}$  for  $\alpha = \overline{\text{id}}$  are very similar. In the following description we focus on  $\Pi^{\text{ip}}$  for  $\alpha = \overline{\text{id}}$ . Removing all conjugates yields the protocol  $\Pi^{\text{ip}}$  for  $\alpha = \text{id}$ . The protocol  $\Pi^{\text{norm}}$  can be obtained by letting the verifier compute the trace of the alleged inner product.

Our approach is based on polynomial identities. That is, for  $\mathbf{w}$ , we define the polynomials  $g(X) = \sum_{i \in [m]} w_i X^i$  and  $\bar{g}(X) = \sum_{i \in [m]} \bar{w}_i X^i$ , and observe that the Laurent polynomial  $L(X) = \sum_{i \in \pm[m]} v_i X^i :=$

$g(X) \cdot \bar{g}(X^{-1})$  has constant coefficient  $\sum_{i \in [m]} w_i \bar{w}_i$ , which is the inner product  $\langle \mathbf{w}, \bar{\mathbf{w}} \rangle_{\mathcal{R}}$ . Also observe that  $v_k = \sum_{i-j=k} v_i \bar{v}_j = \bar{\text{id}} \left( \sum_{i-j=k} \bar{v}_i v_j \right) = \bar{\text{id}} \left( \sum_{j-i=k} \bar{v}_j v_i \right) = \bar{v}_{-k}$  where  $v_k := 0$  if  $|k| \geq m$ . We exploit this symmetry to commit to  $L(X)$  by committing only to  $(v_0, \dots, v_{m-1})$ . Setting  $h(X) = \sum_{i \in [m]} v_i X^i$  and  $\bar{h}(X) = \sum_{i \in [m]} \bar{v}_i X^i$ , we see that  $L(X) = h(X) + \bar{h}(X^{-1}) - v_0$ . We use this equality to prove the polynomial identity  $L(X) = g(X)\bar{g}(X^{-1})$  by evaluating  $g, \bar{g}, h, \bar{h}$  at a random point  $\xi \leftarrow_{\$} \mathcal{C}_{\mathcal{R}_q}$  (and checking if  $v_0 = t$ ).

However, A problem with the soundness occurs if the approach is used naively: The terms  $v_i$  have norm bounded by  $\beta^2$ , so  $\|\mathbf{v}\|$  which may be beyond the threshold for which the commitment is binding.

A natural approach is to run  $\Pi^{b\text{-decomp}}$  to counteract this problem. However, doing so modularly runs into problems and comes at the cost of a suboptimal relaxed knowledge guarantee. We can tighten our analysis if we treat the composition with  $\Pi^{b\text{-decomp}}$  as *within* the protocol  $\Pi^{\text{ip}}$ , i.e. we immediately send the decomposed (and binding) commitments. The reason is a technical artefact of relaxed knowledge soundness and reductions of knowledge: Relaxed soundness in  $\Pi^{b\text{-decomp}}$  incurs a large factor of norm growth, however, in  $\Pi^{\text{ip}}$ , we do not care about the auxiliary commitment to  $\mathbf{v}$  (which norms up to  $\approx \sqrt{m}\beta^2$ ). Thus, we can argue directly for the decomposition of  $\mathbf{v}$  into smaller  $\mathbf{v}_i$  of norm at most  $\beta$ , which are binding. This avoids the need for recovering recover  $\mathbf{v}$  via relaxed knowledge for  $\Pi^{b\text{-decomp}}$ , which significantly improves the parameters. This optimised protocol is presented in Fig. 8

**Lemma 7 (Security of  $\Pi^{\text{norm}}$  and  $\Pi^{\text{ip}}$ ).** *The protocol  $\Pi^{\text{norm}}$  (resp.  $\Pi^{\text{ip}}$ ) is a reduction of knowledge from  $\Xi_{n_{\text{out}}}^{\text{norm}}$  (resp.  $\Xi_{n_{\text{out}}}^{\text{ip}}$ ) to  $(\Xi_{n_{\text{out}+3}^{\text{lin}}})^{\ell+1}$ , where  $b_{\text{ip}} \leq 2\beta/(\sqrt{m}\varphi^{3/2})$  into  $\ell \geq \log_{b_{\text{ip}}}(2\beta^2 + 1)$  for  $b_{\text{ip}}, \ell \in \mathbb{N}$ . It is perfectly correct. It is relaxed knowledge sound from  $\Xi_{n_{\text{out}}}^{\text{norm}\vee\text{sis}}$  (resp.  $\Xi_{n_{\text{out}}}^{\text{ip}\vee\text{sis}}$ ) to  $(\Xi_{n_{\text{out}+3}^{\text{lin}\vee\text{sis}}})^{\ell}$  if  $2\beta \leq \beta^{\text{sis}}$ . Extraction requires two uniformly distributed transcripts (Section 3.2) and has knowledge error  $\kappa \leq \frac{2m}{|\mathcal{C}_{\mathcal{R}_q}|}$ .*

*Proof.* To show that  $\|\mathbf{v}_i\|_2 \leq \beta$ , it suffices to recall the bound for  $\Pi^{b\text{-decomp}}$  from Lemma 3. Indeed, we set  $b_{\text{ip}}$  such that  $\beta \geq \frac{1}{2}\sqrt{m}\varphi^{3/2}b_{\text{ip}}$  holds by definition. That the linear relations hold was already explained above. Thus, the reduction of knowledge is perfectly complete.

Now, we show relaxed knowledge soundness. First, observe that, as argued above, given fixed  $\mathbf{v}$  which defines  $h(X) = \sum_{i=0}^{m-1} v_i X^i$ , then the probability that

$$L(X) \neq h(X) + \bar{h}(X^{-1}) - v_0 \quad \text{but} \quad L(\xi) \neq h(\xi) + \bar{h}(\xi^{-1}) - v_0 \quad (6)$$

holds, is bounded by  $\frac{2m-1}{|\mathcal{C}_{\mathcal{R}_q}|} \leq \frac{2m}{|\mathcal{C}_{\mathcal{R}_q}|} = \kappa$ , where  $\xi \leftarrow_{\$} \mathcal{C}_{\mathcal{R}_q}$ . (By the lemma of Schwartz-Zippel, analogous to

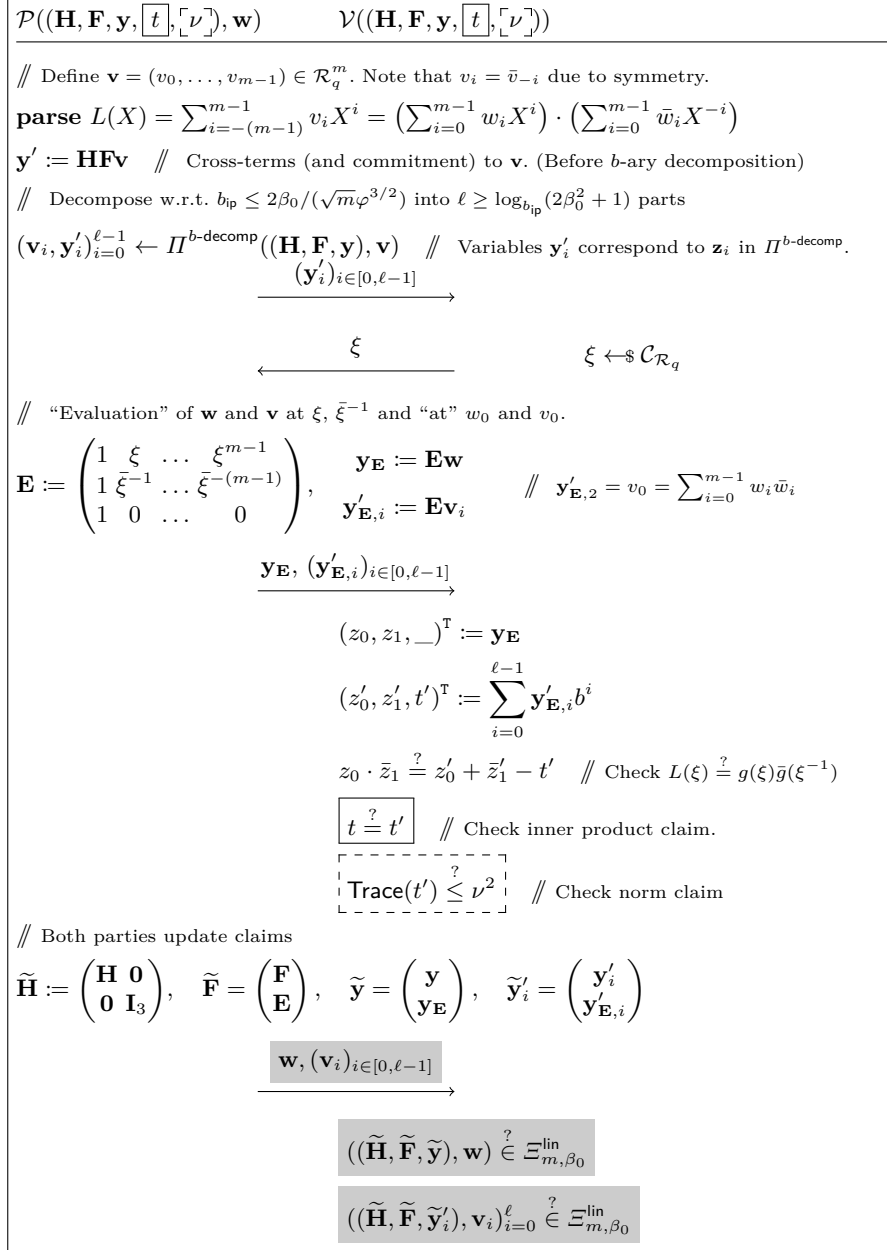
Lemma 6.) However, this is a soundness argument for *fixed* polynomials, but the vector  $\mathbf{v} = \sum_{i=0}^{\ell-1} b_{\text{ip}}^i \mathbf{v}_i$  is only determined in the last step. Now, we argue analogous to Lemma 6 to turn soundness into knowledge soundness: Let  $\mathbf{v}^{(0)}$  with  $\mathbf{HF}\mathbf{v}^{(0)} = \tilde{\mathbf{y}}^{(0)}$  denote the (first) preimage (from the two accepting transcripts), where  $\tilde{\mathbf{y}}^{(0)} = \sum_{i=0}^{\ell} b_{\text{ip}}^i \tilde{\mathbf{y}}_i^{(0)}$  denotes the non-decomposed linear claim for  $\mathbf{v}$ . Observe that only fraction  $\kappa$  of challenges can satisfy (6) for any fixed  $\mathbf{v}^{(0)}$ . Thus, if the adversary succeeds with probability  $\epsilon$ , then with probability  $\epsilon - \kappa$ , the 2-transcript extractor gives transcripts with<sup>16</sup>  $\mathbf{v}_i^{(0)} \neq \mathbf{v}_i^{(1)}$  for some  $i^*$ . Let  $\mathbf{u} = \mathbf{v}_{i^*}^{(0)} - \mathbf{v}_{i^*}^{(1)}$ . Then we have  $\mathbf{HF}\mathbf{u} = 0$ , as  $\mathbf{HF}\mathbf{v}_{i^*}^{(j)} = \mathbf{y}_i^j$  holds for both  $j = 0, 1$ . This yields an OR-branch witness for  $\Xi^{\text{ip}\vee\text{sis}}$  with norm at most  $2\beta \leq \beta^{\text{sis}}$ .

To conclude the proof, we note that we do not care about the auxiliary  $\mathbf{v}_i$ , and they do not play a role in the bound  $\beta$  (aside from the OR-branch), which is for  $\mathbf{w}$  (which is merely passed through by  $\Pi^{\text{norm}}$  and  $\Pi^{\text{ip}}$ ).  $\square$

We recall the technical peculiarity in the construction and proof, that since the auxiliary commitment to  $\mathbf{v}$  in  $\Pi^{\text{norm}}$  and  $\Pi^{\text{ip}}$  is not extracted through  $\Pi^{b\text{-decomp}}$ , but we directly committed to the decomposed  $\mathbf{v}_i$ . For this reason, relaxed soundness incurs almost no norm growth. If we had considered the intermediate relation for  $\mathbf{v}$ , then technically  $\mathbf{v}$  becomes part of the witness which the reduction of knowledge for  $\Pi^{b\text{-decomp}}$  must then extract. This blows up the norm bound (by  $\beta$ ).

**The full protocols  $\Pi^{\text{norm}+}$  and  $\Pi^{\text{ip}+}$ .** The full protocol  $\Pi^{\text{norm}+}$  (resp.  $\Pi^{\text{ip}+}$ ) simply reduces the  $\ell + 1$  statements after  $\Pi^{\text{norm}}$  (resp.  $\Pi^{\text{ip}}$ ) back to a single one with the minimal number of rows, by applying  $\Pi^{\text{batch}}$  and  $\Pi^{\text{fold}}$ . That is,  $\Pi^{\text{norm}+}$  (and analogously  $\Pi^{\text{ip}+}$ ) work as follows:

<sup>16</sup>We exploit that the challenges are *uniformly* distributed (conditioned on accepting).



**Fig. 8.** Protocols  $\boxed{\Pi^{\text{norm}}}$  for  $\Xi^{\text{norm}}$  and  $\boxed{\Pi^{\text{ip}}}$  for  $\Xi^{\text{ip}}$  where  $\alpha = \bar{\text{id}}$ . The case  $\alpha = \text{id}$  can be obtained by removing all conjugates.  $\Pi^{\text{norm}}$  and  $\Pi^{\text{ip}}$  send the marked parts only as a *proof* (but not *reduction*) of knowledge.

- (1)  $\Xi_{n^{\text{out}},\beta}^{\text{norm}} \xrightarrow{II^{\text{norm}}} (\Xi_{n^{\text{out}}+3,\beta}^{\text{lin}})^{\ell+1}$  (Reduce to claims in  $(\Xi^{\text{lin}})^{\ell+1}$ .)
- (2)  $(\Xi_{n^{\text{out}}+3,\beta}^{\text{lin}})^{\ell+1} \xrightarrow{II^{\text{batch}}} (\Xi_{\bar{n}+1,\beta}^{\text{lin}})^{\ell+1}$  (Reduce number of rows.)
- (3)  $(\Xi_{\bar{n}+1,\beta}^{\text{lin}})^{\ell+1} \xrightarrow{II^{\text{fold}}} \Xi_{\bar{n}+1,(\ell+1)\gamma\beta}^{\text{lin}}$  (Reduce to one claim with  $\beta' = (\ell+1)\gamma\beta$ .)

In words, first  $II^{\text{norm}+}$  (resp.  $II^{\text{ip}+}$ ) runs  $II^{\text{norm}}$  (resp.  $II^{\text{ip}}$ ) to obtain the commitments to  $\mathbf{v}_i$ . Next it runs  $II^{\text{batch}}$  to reduce to  $\bar{n}+1$  rows again (i.e. to  $\underline{n} = 1$ ). Finally, the claims are folded into one by using  $II^{\text{fold}}$ . The full norm protocol  $II^{\text{norm}}$  and its security guarantees are analogous.

*Remark 6.* When applying  $II^{\text{batch}}$  to the product relation  $(\Xi_{n^{\text{out}}+3,\beta}^{\text{lin}})^{\ell+1}$ , it is crucial a single challenge is reused among all instances. This ensures that if all  $(\mathbf{H}_i, \mathbf{F}_i)$  were identical before batching, they still are identical after batching.

**Theorem 6 (Security of  $II^{\text{norm}+}$  and  $II^{\text{ip}+}$ ).** *The protocol  $II^{\text{norm}+}$  (resp.  $II^{\text{ip}+}$ ) is a reduction of knowledge  $\Xi_{n^{\text{out}},\beta_0}^{\text{norm}}$  (resp.  $\Xi_{n^{\text{out}},\beta_0}^{\text{ip}}$ ) to  $(\Xi_{\bar{n}+1,\beta_3}^{\text{lin}})^{\ell+1}$ , where  $\ell \geq \log_{b_{\text{ip}}}(2\beta^2 + 1)$  for  $b_{\text{ip}} \leq 2\beta/(\sqrt{m}\varphi^{3/2})$  and  $\beta_3 = (\ell+1)\beta_0$  for  $b_{\text{ip}}, \ell \in \mathbb{N}$ . It is perfectly correct. It is relaxed knowledge sound from  $\Xi_{n^{\text{out}},\beta'_0}^{\text{norm}\vee\text{sis}}$  (resp.  $\Xi_{n^{\text{out}},\beta'_0}^{\text{ip}\vee\text{sis}}$ ) to  $\Xi_{\bar{n}+1,\beta'_3}^{\text{lin}\vee\text{sis}}$  if  $\beta'_0 \leq \beta^{\text{sis}}$ , where  $\beta'_0 = 4\theta\beta'_3$ . The knowledge error is  $\frac{\ell+1}{|\mathcal{C}_{\mathcal{R}}|} + \frac{4(\ell+1)+2m}{|\mathcal{C}_{\mathcal{R}_q}|}$ .*

*Proof.* For correctness, let  $\beta_i$  denote the norm bound after the  $i$ -th protocol was applied. Then we get We have

- $\beta_1 = \beta_0$  by Lemma 7.
- $\beta_2 = \beta_1$  by Lemma 6.
- $\beta_3 = \gamma(\ell+1)\beta_2$  by Lemma 5.

Overall, we find

$$\beta_3 = \gamma(\ell+1)\beta_0.$$

For knowledge soundness, we work in the opposite direction. Starting out from a witness in  $\Xi_{m/d,\beta'_3}^{\text{lin}}$ , we follow the bounds back to the start. We have

- $\beta'_2 = 2\theta\beta'_3$  by Lemma 5.
- $\beta'_1 = 2\beta'_2$  by Lemma 6.
- $\beta'_0 = 2\beta'_1$  by Lemma 7.

This yields overall

$$\beta'_0 = 8\theta\beta'_3.$$

The knowledge error for  $II^{\text{ip}}$  is  $\frac{2m}{|\mathcal{C}_{\mathcal{R}_q}|}$ , For  $II^{\text{batch}}$  it is  $\frac{4}{|\mathcal{C}_{\mathcal{R}_q}|}$  per instance (as 4 rows are batched, namely  $\underline{n} + 3$ ), so  $\frac{4(\ell+1)}{|\mathcal{C}_{\mathcal{R}_q}|}$  overall. For  $II^{\text{fold}}$  it is  $\frac{\ell+1}{|\mathcal{C}_{\mathcal{R}}|}$  Taken together, we arrive at the claimed knowledge error. Note that the case of a vSIS break is handled implicitly in the knowledge reductions.  $\square$

**Norm Checks to Upgrade Relaxed Soundness.** Currently, our norm check is defined for the relation  $\Xi^{\text{norm}}$ , which inherits the parameter  $\beta$  from  $\Xi^{\text{lin}}$ , and also contains an explicit  $\nu$  as a norm *statement*. For convenience, we now define the  $II^{\text{norm}\beta}$  protocol, which is a reduction of knowledge from  $\Xi_{\beta}^{\text{lin}}$  to  $\Xi_{\beta}^{\text{lin}}$ , that works as follows:  $II^{\text{norm}\beta}$  runs  $II^{\text{norm}}$  on (implicit) input  $\nu = \beta$  for both parties.

Analogous to  $II^{\text{norm}+}$ , we define  $II^{\text{norm}\beta+} : \Xi_{\beta}^{\text{lin}} \rightarrow \Xi_{\beta}^{\text{lin}}$ . Clearly,  $II^{\text{norm}\beta}$  (resp. and  $II^{\text{norm}\beta+}$ ) directly inherits all correctness and security guarantees of  $II^{\text{norm}}$  (resp.  $II^{\text{norm}+}$ ). We introduce  $II^{\text{norm}\beta}$  and  $II^{\text{norm}\beta+}$  solely for compositional reasons: They start and end with a claim(s) in  $\Xi_{\beta}^{\text{lin}}$ . Let us stress that the  $II^{\text{norm}\beta}$  protocol *upgrades* the (previously relaxed) bound on the norm of  $\mathbf{w}$  to  $\|\mathbf{w}\| \leq \beta$ . We capture this in following corollary.

**Corollary 1.** *Adopt the setting of Theorem 6. Then  $II^{\text{norm}\beta+}$  is relaxed knowledge sound from  $\Xi_{\beta}^{\text{lin}\vee\text{sis}}$  to  $\Xi_{\beta'}^{\text{lin}\vee\text{sis}}$  if  $2\beta' \leq \beta^{\text{sis}}$ , with the same knowledge error as Theorem 6.*

*Proof.* Unless the OR-branch occurs in  $\Xi^{\text{lin}\vee\text{sis}}$ , the norm of the starting witness  $\mathbf{w}$  is *proven* to be at most  $\nu = \beta$  by  $II^{\text{norm}}$ . Thus, the extracted witness has norm at most  $\beta$  or it a witness for the OR-branch.

## 5.8 $\Pi^{\text{sfm}}$ : Split-and-Fold with Norm Checks

We describe our split-and-fold protocol with intermediate norm check  $\Pi^{\text{sfm}}$ . It first runs the norm check  $\Pi^{\text{norm}^+}$  to upgrade the relaxed norm bound to a strict one. Then it splits-and-folds to reduce the witness size. If parameters are set correctly, then  $\Pi^{\text{split}\&\text{fold}}: \Xi_{m,\beta}^{\text{lin}} \rightarrow \Xi_{m/d,\beta}^{\text{lin}}$  is reduction of knowledge with relaxed knowledge soundness  $\Xi_{m,\beta}^{\text{lin}\vee\text{sis}} \rightarrow \Xi_{m/d,\beta'}^{\text{lin}\vee\text{sis}}$ . Crucially, the bound  $\beta$  is guaranteed exactly after extraction (unless the OR-branch, i.e. a vSIS break is extracted).

- (1)  $\Xi_{\beta_0}^{\text{lin}} \xrightarrow{\Pi^{\text{norm}^+}} \Xi_{\beta_1}^{\text{lin}}$ : Run a norm check for  $\beta_0$ .
- (2)  $\Xi_{m,\beta_1}^{\text{lin}} \xrightarrow{\Pi^{\text{split}\&\text{fold}}} \Xi_{m/d,\beta_2}^{\text{lin}}$ : Run the split-and-fold to reduce witness size.

**Theorem 7.** *Let*

- Let  $b_{\text{ip}} \leq 2\beta/(\sqrt{m}\varphi^{3/2})$  and  $\ell_0 \geq \log_{b_{\text{ip}}}(2\beta^2 + 1)$ .
- Let  $b_1 \geq 1$  be arbitrary and  $\ell_1 = \lceil \log_{b_1}(2\beta_0 + 1) \rceil$ .
- Let  $\beta_1 = \gamma(\ell_0 + 1)\beta_0$ .
- Let  $\beta_2 = \gamma\ell_1 d \frac{1}{2} \sqrt{m}\varphi^{3/2} b_1$ .

Then  $\Pi^{\text{sfm}}$  is perfectly correct from  $\Xi_{m,\beta_0}^{\text{lin}}$  to  $\Xi_{m/d,\beta_2}^{\text{lin}}$ .

If  $\beta'_0 \leq \beta^{\text{sis}}$ , where

- $\beta'_1 = 4\sqrt{d}\theta\beta_0\beta'_2$
- $\beta'_0 = 8\theta\beta'_1$

then  $\Pi^{\text{sfm}}$  is relaxed knowledge sound from  $\Xi_{m,\beta_0}^{\text{lin}\vee\text{sis}}$  to  $\Xi_{m/d,\beta'_2}^{\text{lin}\vee\text{sis}}$ . It has knowledge error  $\frac{2\ell+1}{|\mathcal{C}_{\mathcal{R}}|} + \frac{3(\ell_0+1)+2m+d-1}{|\mathcal{C}_{\mathcal{R}_q}|}$

*Proof.* Correctness follows since  $\beta_1$  is the correctness assertion in Corollary 1 (or rather Theorem 6), and  $\beta_2$  is the correctness assertion in Theorem 5. Similarly, the relaxed knowledge soundness follows by plugging in the respective relaxed knowledge soundness claims. Finally, the knowledge error is the sum of the knowledge errors.  $\square$

$\Pi$	$\beta_1$	$\beta'_0$	$\kappa$	#tr	Other
$\Pi^{\text{b-decomp}}$	$\frac{1}{2}\sqrt{m}\varphi^{3/2}b$	$2\beta_0\beta'_1$	0	1	$\ell = \lceil \log_b(2\beta_0 + 1) \rceil$
$\Pi^{\text{split}}$	$\beta_0$	$\sqrt{d}\beta'_1$	$(d-1)/ \mathcal{C}_{\mathcal{R}} $	$d$	Need $2\beta'_0 \leq \beta^{\text{sis}}$
$\Pi^{\text{fold}}$	$\gamma\ell\beta_0$	$2\theta\beta'_1$	$\ell/ \mathcal{C}_{\mathcal{R}} $	$\ell + 1$	$\ell$ from $(\Xi^{\text{lin}})^\ell$
$\Pi^{\text{batch}}$	$\beta_0$	$\beta'_1$	$\underline{n}/ \mathcal{C}_{\mathcal{R}_q} $	2	Need $2\beta'_0 \leq \beta^{\text{sis}}$
$\Pi^{\text{norm}}/\Pi^{\text{ip}}$	$\beta_0$	$\beta'_1$	$2m/ \mathcal{C}_{\mathcal{R}_q} $	2	
$\Pi^{\text{split}\&\text{fold}}$	$\frac{1}{2}\gamma\ell d\sqrt{m}\varphi^{3/2}b\beta_0$	$4\sqrt{d}\theta\beta_0\beta'_1$	$\ell/ \mathcal{C}_{\mathcal{R}}  + (d-1)/ \mathcal{C}_{\mathcal{R}_q} $	$d(\ell + 1)$	$\ell = \ell_{\Pi^{\text{b-decomp}}}$
$\Pi^{\text{norm}^+}/\Pi^{\text{ip}^+}$	$\gamma(\ell + 1)\beta_0$	$8\theta\beta'_1$	$\frac{\ell+1}{ \mathcal{C}_{\mathcal{R}} } + \frac{4(\ell+1)+2m}{ \mathcal{C}_{\mathcal{R}_q} }$	$4(\ell + 2)$	See caption.

**Table 1.** Parameters of protocols. Expressed as  $\beta_1 = f(\beta_0)$  for correctness,  $\beta'_0 = g(\beta'_1)$  and  $\beta'_0 \leq \beta^{\text{sis}}$  for relaxed soundness, knowledge error  $\kappa$ , number of transcripts to extract, other variables. For  $\Pi^{\text{ip}^+}$ , we have  $\ell = \log_{b_{\text{ip}}}(2\beta_0^2 + 1)$  for  $b_{\text{ip}} \leq 2\beta/(\sqrt{m}\varphi^{3/2})$ .

## 5.9 Asymptotic Communication Complexity

We now compute the proof size for the split-and-fold with norm checks protocol in Section 5.8. Having non-interactive proof systems in mind, we only count prover messages. Let  $m = d^\mu$  where  $d = O(1)$  and  $\mu = O(\log m)$ . The other parameters are chosen according to Table 1. As shown in Section 4, we can pick  $f = \text{poly}(\lambda)$  and a subtractive set  $\mathcal{C}_{\mathcal{R}}$  such that  $\gamma = O(1)$  and  $\theta = \text{poly}(\lambda)$ .

Recall the prover executes two sub-protocols:



- (1)  $\Xi_{\beta_0}^{\text{lin}} \xrightarrow{\Pi^{\text{norm}^+}} \Xi_{\beta_1}^{\text{lin}}$ : Run a norm check for  $\beta_0$ .
- (2)  $\Xi_{m,\beta_1}^{\text{lin}} \xrightarrow{\Pi^{\text{split}\&\text{fold}}} \Xi_{m/d,\beta_2}^{\text{lin}}$ : Run the split-and-fold to reduce witness size.

First, we turn to the  $\Pi^{\text{norm}^+}$  protocol in Figure 8. To prove relation  $\Xi_{\beta_0}^{\text{lin}}$ , the prover starts with

$$\Xi_{n^{\text{out}},\beta}^{\text{norm}} \xrightarrow{\Pi^{\text{norm}}} (\Xi_{n^{\text{out}}+3,\beta}^{\text{lin}})^{\ell+1}$$

where it sends all  $\mathbf{y}'_i$  which in total have size  $\ell_0 n^{\text{out}}$  ring elements. After receiving the challenge  $\xi$ , the prover outputs  $\mathbf{y}_E$  and  $\mathbf{y}'_{E,i}$  of size 3 each – thus in total  $3(\ell_0 + 1)$  ring elements. Then the prover runs

$$(\Xi_{n^{\text{out}}+3,\beta}^{\text{lin}})^{\ell+1} \xrightarrow{\Pi^{\text{batch}}} (\Xi_{\bar{n}+1,\beta}^{\text{lin}})^{\ell+1} \xrightarrow{\Pi^{\text{fold}}} \Xi_{\bar{n}+1,(\ell+1)\gamma\beta}^{\text{lin}}$$

where neither  $\Pi^{\text{batch}}$  nor  $\Pi^{\text{fold}}$  incur any no communication from the prover.

Next, we now move on to the  $\Pi^{\text{split}\&\text{fold}}$  protocol. The prover wants to give a proof for relation  $\Xi_{m,\beta_1}^{\text{lin}}$ . The prover starts by running the  $\Pi^{b\text{-decomp}}$  protocol

$$\Xi_{m,\beta_1}^{\text{lin}} \xrightarrow{\Pi^{b\text{-decomp}}} (\Xi_{m,b_1}^{\text{lin}})^{\ell_1}$$

where the prover sends  $\ell_1$  vectors  $(\mathbf{z}_k)$  of size  $n^{\text{out}}$  elements in  $\mathcal{R}_q$ . Next, it runs  $\Pi^{\text{split}}$

$$(\Xi_{m,b_1}^{\text{lin}})^{\ell_1} \xrightarrow{\Pi^{\text{split}}} (\Xi_{m/d,b_1}^{\text{lin}})^{\ell_1 d}$$

and outputs (using the optimization in Remark 5): (i)  $d - 1$  intermediate “top-part” evaluations  $\bar{\mathbf{y}}_j$  of size  $\bar{n}$  elements in  $\mathcal{R}_q$  and (ii)  $d^2 - 1$  “bottom-part” evaluations  $\underline{\mathbf{y}}_{i,j}$  of size  $\underline{n}$  elements in  $\mathcal{R}_q$ . Finally the prover executes the  $\Pi^{\text{fold}}$  protocol

$$(\Xi_{m/d,b_1}^{\text{lin}})^{\ell_1 d} \xrightarrow{\Pi^{\text{fold}}} \Xi_{m/d,\beta_2}^{\text{lin}}$$

where no prover message is sent. All in all, in each of  $\mu = O(\log m)$  iterations of split-and-fold with norm checks, the prover sends

$$(\ell_0 n^{\text{out}} + 3(\ell_0 + 1)) + (\ell_1 n^{\text{out}} + (d - 1)\bar{n} + (d^2 - 1)\underline{n})$$

elements in  $\mathcal{R}_q$ .

**Simple example: vSIS opening proof.** When proving knowledge of a vSIS commitment opening, we can set  $n, n^{\text{out}} \in O(1)$ . Due to the polynomial challenge space for subtractive sets, we need to repeat  $O(\lambda/\log \lambda)$  times to ensure  $\approx 2^{-\lambda}$  soundness error. Finally, we can pick  $b_1$  such that  $\ell_1 = O(1)$ . In total, the proof size in the number of ring elements is simply bounded by  $O\left(\lambda \frac{\log m}{\log \lambda}\right)$ .

Next, we turn to setting the bound required for the (v)SIS problem to be hard. From Theorem 7 we need to set  $\beta'_2 = \beta_2$  and

$$\beta^{\text{sis}} = \beta'_0 = 32\sqrt{d}\theta^2 \beta_0 \beta'_2 = 16\gamma \ell_1 d^{3/2} \sqrt{m} \varphi^{3/2} \theta^2 \beta_0 b_1 = \text{poly}(m, \lambda).$$

The next step is to estimate an asymptotic size of a ring element in  $\mathcal{R}_q$ .

**Hardness of SIS.** To measure hardness of vSIS, we heuristically assume that it is as hard as the plain SIS problem for the dimension  $\varphi = \varphi(\mathbf{f})$ . To measure hardness of SIS, we first translate the canonical norm  $\|\sigma(\cdot)\|_2$  into the Euclidean norm  $\|\psi(\cdot)\|_2$ , and then follow the heuristic methodology from [MR09]. That is, let  $b = O(\lambda)$  be the block size of the BKZ algorithm to find a short vector in the corresponding  $q$ -ary lattice for SIS (cf. [BDGL15]). Define the root Hermite factor as  $\delta_{\text{rhf}} = \left(\frac{b(\pi b)^{1/b}}{2\pi e}\right)^{1/(2(b-1))}$ . Then, SIS with matrix dimensions  $\varphi \times \varphi m$  and Euclidean norm  $\beta^* = \beta^{\text{sis}} \cdot \text{poly}(\lambda)$  is hard when  $\beta^* < \min\left(2^{2\sqrt{\varphi \log q \log \delta_{\text{rhf}}}}, q\right)$ . By rearranging, we get that  $\varphi \log q > \log^2 \beta^* / 4 \log \delta_{\text{rhf}}$ . Note that

$$\log \delta_{\text{rhf}} = \frac{1}{2(b-1)} \log \left(\frac{b(\pi b)^{1/b}}{2\pi e}\right) = \Theta\left(\frac{\log b}{b}\right) = \Theta\left(\frac{\log \lambda}{\lambda}\right).$$

Finally, using the fact that  $\beta^* = \text{poly}(m, \lambda)$ , size of a single  $\mathcal{R}_q$  element is asymptotically

$$\Omega\left(\frac{\lambda \cdot (\log m + \log \lambda)^2}{\log \lambda}\right) = \Omega\left(\lambda \cdot \left(\frac{\log^2 m}{\log \lambda} + \log \lambda\right)\right) \text{ bits.}$$

Therefore, we deduce that the total proof size in bits is  $O(\log^3 m \cdot \lambda^2 / \log^2 \lambda)$ .

## 6 Packed $\mathbb{Z}$ -Inner Products via Twisted Trace Maps

We propose an abstract framework based on “twisted trace maps” that reduces  $\mathbb{Z}$ -inner products to  $\mathcal{R}$ -inner products over various choices of  $\mathcal{R}$ . In a nutshell, for a fixed choice of  $\mathcal{R}$ , we would like to construct a twisted trace map  $\tau : \mathcal{R} \rightarrow \mathbb{Z}$  of the form shown below, where  $N \in \mathbb{N}$  is some normalisation factor and  $\alpha \in \mathcal{R}$  is called a “twist” element, such that the following diagram commutes:

$$\begin{array}{ccc} \mathbb{Z}^\delta \times \mathbb{Z}^\delta & \xrightarrow{\langle \cdot, \cdot \rangle_{\mathbb{Z}}} & \mathbb{Z} \\ \psi^{-1}(\cdot) \times \overline{\psi^{-1}(\cdot)} \downarrow & & \uparrow \tau \\ \mathcal{R} \times \mathcal{R} & \xrightarrow{\cdot_{\mathcal{R}}} & \mathcal{R} \end{array} \quad \text{where} \quad \tau : z \mapsto \frac{1}{N} \cdot \text{Trace}(\alpha \cdot z).$$

**Definition 5 (Inner-Product Embedding).** Let  $\mathcal{R} \subset \mathcal{O}_{\mathcal{K}}$  be a subring identified by a  $\mathbb{Z}$ -basis  $\mathbf{b} \in \mathcal{R}^\delta$  of  $\delta$  elements. We say that a tuple  $\tau$  is an inner-product embedding over  $\mathcal{R}$  if  $\tau : \mathcal{R} \rightarrow \mathbb{Z}^\delta$  is a  $\mathbb{Z}$ -linear map and, for any  $\mathbf{x}, \mathbf{y} \in \mathbb{Z}^\delta$ , it holds that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \tau \left( \psi_{\mathbf{b}}^{-1}(\mathbf{x}) \cdot_{\mathcal{R}} \overline{\psi_{\mathbf{b}}^{-1}(\mathbf{y})} \right)$ .

### 6.1 Power-of-Two Cyclotomics via Constant Term

As a simple concrete example, we recall a well-known folklore technique for computing the inner product over the coefficient embeddings of power-of-two cyclotomics [LNP22].

**Theorem 8.** Let  $\mathcal{R} = \mathbb{Z}[\zeta_{\mathfrak{f}}]$  with a conductor  $\mathfrak{f} = 2^k$  for some  $k \in \mathbb{N}$ ,  $\delta = \varphi = \varphi(\mathfrak{f}) = \mathfrak{f}/2$ ,  $\tau(\cdot) = \text{ct}(\cdot) = (\psi(\cdot))_0$ , where  $\psi$  denotes the coefficient embedding and  $\text{ct}(\cdot)$  is the constant term of the coefficient embedding. Then  $\tau$  is an inner-product embedding over  $\mathcal{R}$ .

*Proof.* Write  $\zeta = \zeta_{\mathfrak{f}}$ . Observe that  $\text{ct}(\zeta^i) = (i \stackrel{?}{=} 0)$  for all  $i \in \pm[\varphi]$ .<sup>17</sup> Consider the vectors  $\mathbf{x} = (x_0, \dots, x_{\varphi-1}), \mathbf{y} = (y_0, \dots, y_{\varphi-1}) \in \mathbb{Z}^\varphi$ . The elements  $x := \psi^{-1}(\mathbf{x})$  and  $\bar{y} := \overline{\psi^{-1}(\mathbf{y})}$  can be expressed as

$$x = \psi^{-1}(\mathbf{x}) = \sum_{i \in [\varphi]} x_i \zeta^i \quad \text{and} \quad \bar{y} = \overline{\psi^{-1}(\mathbf{y})} = \sum_{i \in [\varphi]} y_i \zeta^{-i}$$

respectively. Note that their product  $z := x \cdot \bar{y}$  satisfies

$$\begin{aligned} z = x \cdot \bar{y} &= \left( \sum_{i \in [\varphi]} x_i \zeta^i \right) \left( \sum_{j \in [\varphi]} y_j \zeta^{-j} \right) = \sum_{i, j \in [\varphi]} x_i y_j \zeta^{i-j} \\ &= \underbrace{\sum_{i \in [\varphi]} x_i y_i}_{\text{ct}(z)} + \sum_{i, j \in [\varphi]: i \neq j} x_i y_j \zeta^{i-j}. \end{aligned}$$

Therefore,  $\tau(\psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{y})}) = \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}}$ , as desired.  $\square$

*Remark 7.* The constant term map  $\text{ct}(x)$  from Theorem 8 can be expressed in terms of the Trace function as  $\tau(x) = \frac{1}{\varphi} \text{Trace}(x)$ , where  $\varphi = \mathfrak{f}/2$  since  $\mathfrak{f}$  is a power of 2, and may be viewed as a twisted trace map  $\tau(x) = \frac{1}{\varphi} \text{Trace}(\alpha \cdot x)$  with  $\alpha = 1$ .

As pointed out in Section 4, power-of-two cyclotomic rings do not admit large subtractive sets, and are therefore ill-suited for certain applications, e.g. instantiating the succinct arguments presented in Section 5. This motivates the search for inner-product embeddings  $\tau$  over other rings.

<sup>17</sup>This is true for power-of-2 cyclotomics since power-of-2 cyclotomic polynomials are of the form  $\Phi_{\mathfrak{f}}(X) = X^\varphi + 1$ . Note that this is false for non-power-of-2 conductors. For example, if  $\mathfrak{f}$  is prime, then  $\zeta^{-1} = \zeta^\varphi = -\sum_{i \in [\varphi]} \zeta^i$  with  $\text{ct}(\zeta^{-1}) = -1$ .

## 6.2 Prime Real Cyclotomics via Twisted Trace

A natural class of rings to search for inner-product embeddings are cyclotomic rings with large prime conductors, since they admit large subtractive sets (cf. Section 4). Although we did not manage to design inner-product embeddings in those rings, we did so for its maximal real subring, adapting a result from lattice code theory [BFOV04, Proposition 1].

**Theorem 9.** *Let  $\mathcal{K} = \mathbb{Q}(\zeta_{\mathfrak{f}})$  where  $\mathfrak{f}$  is prime and  $\mathcal{R} = \mathbb{Z}[\zeta_{\mathfrak{f}} + \zeta_{\mathfrak{f}}^{-1}]$  be identified by the  $\mathbb{Z}$ -basis  $\mathbf{b}^+ = \left\{ \sum_{i=[j+1]}^{\varphi/2-i} (\zeta^{\varphi/2-i} + \zeta^{-(\varphi/2-i)}) \right\}_{j \in [\varphi/2]}$ . For  $z \in \mathcal{R}$ , let  $\tau(z) = \frac{1}{2\mathfrak{f}} \text{Trace}(\alpha z)$  be a twisted trace map for the twist element  $\alpha = t \cdot \bar{t}$  where  $t = \zeta^{-\varphi/2} - \zeta^{\varphi/2}$ . Then  $\tau$  is an inner product embedding over  $\mathcal{R}$ .*

*Proof.* Since  $\mathfrak{f}$  is prime, we have  $\varphi = \mathfrak{f} - 1$  and  $\delta = \varphi/2 = (\mathfrak{f} - 1)/2$ . In the following, write  $\text{Trace} = \text{Trace}$ . Recall that  $\text{Trace}(1) = \sum_{j \in [\varphi]} 1 = \varphi = \mathfrak{f} - 1$ . Furthermore, for  $i \in \mathbb{Z}_{\mathfrak{f}}^{\times}$ , we have  $\text{Trace}(\zeta^i) = \sum_{j \in \mathbb{Z}_{\mathfrak{f}}^{\times}} \zeta^{ij} = \sum_{j \in \mathbb{Z}_{\mathfrak{f}}^{\times}} \zeta^j = -1$ .

As a starting point, we consider the following sequence.

$$\mathbf{b}^- = (b_i^-)_{i \in [\varphi/2]} = (\zeta^{i+1} - \zeta^{-(i+1)})_{i \in [\varphi/2]}.$$

Note that the sequence  $\mathbf{b}^-$  is trace-orthogonal. Namely, the following function acts as Kronecker delta.

$$\frac{1}{2\mathfrak{f}} \cdot \text{Trace}(\overline{b_i^-} \cdot b_j^-) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

To see this, let  $i' = i + 1$  and  $j' = j + 1$  and consider the trace of the expression below.

$$\begin{aligned} \overline{b_i^-} \cdot b_j^- &= \overline{(\zeta^{i'} - \zeta^{-i'})} \cdot (\zeta^{j'} - \zeta^{-j'}) \\ &= (\zeta^{-i'} - \zeta^{i'}) \cdot (\zeta^{j'} - \zeta^{-j'}) = (\zeta^{-i'+j'} - \zeta^{i'+j'} - \zeta^{-i'-j'} + \zeta^{i'-j'}), \\ \text{Trace}(\overline{b_i^-} \cdot b_j^-) &= \text{Trace}(\zeta^{-i'+j'} - \zeta^{i'+j'} - \zeta^{-i'-j'} + \zeta^{i'-j'}) \\ &= \text{Trace}(\zeta^{-(i'-j')}) - \text{Trace}(\zeta^{i'+j'}) - \text{Trace}(\zeta^{-(i'+j')}) + \text{Trace}(\zeta^{i'-j'}). \end{aligned}$$

Since  $i, j \in [\varphi/2]$ , we have  $2 \leq i' + j' \leq \varphi$ , meaning that  $i' + j' \in \mathbb{Z}_{\mathfrak{f}}^{\times}$  and  $-(i' + j') \in \mathbb{Z}_{\mathfrak{f}}^{\times}$ . Furthermore, if  $i \neq j$ , then  $i' - j' \in \pm[\varphi/2] \setminus \{0\}$ , hence  $i' - j' \in \mathbb{Z}_{\mathfrak{f}}^{\times}$  and  $-(i' - j') \in \mathbb{Z}_{\mathfrak{f}}^{\times}$ . Therefore, we conclude that

$$\begin{aligned} (i = j) &\implies \text{Trace}(\overline{b_i^-} \cdot b_j^-) = 2\mathfrak{f}, \\ (i \neq j) &\implies \text{Trace}(\overline{b_i^-} \cdot b_j^-) = 0. \end{aligned}$$

Although  $\mathbf{b}^-$  is trace-orthogonal, it does not constitute a basis of any ring. It does, however, match in cardinality the degree of the maximal real subring  $\mathcal{R}$ , to which our attention now turns. Consider the “suffix-sum” basis for the maximal real subring as in the theorem statement:

$$\mathbf{b}^+ = (b_j^+)_{j \in [\varphi/2]} = \left\{ \sum_{i \in [j+1]}^{\varphi/2-i} (\zeta^{\varphi/2-i} + \zeta^{-(\varphi/2-i)}) \right\}_{j \in [\varphi/2]}$$

From the identity  $-\varphi/2 = \mathfrak{f} - \varphi/2 = \varphi + 1 - \varphi/2 = \varphi/2 + 1 \pmod{\mathfrak{f}}$ , for each  $j \in [\varphi/2]$ , we observe that

$$b_j^- = b_j^+ \cdot \underbrace{(\zeta^{-\varphi/2} - \zeta^{\varphi/2})}_t$$

since

$$b_j^+ \cdot (\zeta^{-\varphi/2} - \zeta^{\varphi/2}) = \sum_{i \in [j+1]}^{\varphi/2-i} (\zeta^{\varphi/2-i} + \zeta^{-(\varphi/2-i)}) \cdot (\zeta^{-\varphi/2} - \zeta^{\varphi/2})$$

$$\begin{aligned}
&= \sum_{i \in [j+1]} (\zeta^{-i} - \zeta^i + \zeta^{-(\varphi-i)} - \zeta^{\varphi-i}) \\
&= \sum_{i \in [j+1]} (\zeta^{-i} - \zeta^{-(i+1)} + \zeta^{i+1} - \zeta^i) \\
&= \sum_{i \in [j+1]} (\zeta^{-i} - \zeta^{-(i+1)}) + \sum_{i \in [j+1]} (\zeta^{i+1} - \zeta^i) \\
&= \zeta^{j+1} - \zeta^{-j+1} = b_j^-.
\end{aligned}$$

Therefore,

$$\frac{1}{2\mathfrak{f}} \cdot \text{Trace}(t \cdot b_i^+ \cdot \overline{t \cdot b_j^+}) = \frac{1}{2\mathfrak{f}} \cdot \text{Trace}(\alpha \cdot b_i^- \cdot \overline{b_j^-}) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

Now, suppose  $x = \psi_{\mathfrak{b}^+}^{-1}(\mathbf{x})$  and  $\bar{y} = \overline{\psi_{\mathfrak{b}^+}^{-1}(\mathbf{y})}$  for some  $\mathbf{x}, \mathbf{y} \in \mathbb{Z}^\delta$ . We have

$$\begin{aligned}
\tau(x \cdot \bar{y}) &= \frac{1}{2\mathfrak{f}} \text{Trace}(\alpha x \cdot \bar{y}) = \sum_{i, j \in [\varphi/2]} x_i y_j \frac{1}{2\mathfrak{f}} \text{Trace}(t \cdot b_i^+ \cdot \bar{t} \cdot \overline{b_j^+}) = \sum_{i \in [\varphi/2]} x_i y_i \\
&= \langle \mathbf{x}, \mathbf{y} \rangle. \quad \square
\end{aligned}$$

The above theorem constructs inner-product embeddings for  $\mathcal{R} = \mathbb{Z}[\zeta_{\mathfrak{f}} + \zeta_{\mathfrak{f}}^{-1}]$  where  $\mathfrak{f}$  is prime. This restricts the choice of  $\mathcal{R}$  quite severely, especially considering that the subtractive set constructed in Section 4 for  $\mathbb{Z}[\zeta_{\mathfrak{f}}]$  or  $\mathbb{Z}[\zeta_{\mathfrak{f}} + \zeta_{\mathfrak{f}}^{-1}]$  for prime  $\mathfrak{f}$  has a large expansion factor bound  $\gamma_S \leq \mathfrak{f}$ .

### 6.3 Tensor of Prime Real Cyclotomics

To allow more fine-grained parameter selection, we extend the result in Sections 6.1 and 6.2 by constructing larger rings using the tensor product, inspired by [BFOV04, Proposition 6]. Concretely, we construct subtractive sets for rings  $\mathcal{R} = \mathcal{O}_{\mathcal{K}_{2^d}} \otimes \mathcal{O}_{\mathcal{K}_{\mathfrak{f}_0}^+} \otimes \dots \otimes \mathcal{O}_{\mathcal{K}_{\mathfrak{f}_{k-1}}^+}$  for distinct odd primes  $\mathfrak{f}_0, \dots, \mathfrak{f}_{k-1}$ . Note that  $\mathcal{R}$  has conductor  $\mathfrak{f} = 2^d \cdot \prod_{i \in [k]} \mathfrak{f}_i$  and degree  $\delta = 2^d \cdot \prod_{i \in [k]} (\mathfrak{f}_i - 1)$ . It is contained in the ring  $\mathcal{O}_{\mathcal{K}_{\mathfrak{f}}}$  which admits a subtractive set  $S$  of size  $\mathfrak{f}/\mathfrak{f}_{\max}$  with expansion factor  $\gamma_S = 1$  (cf. Section 4).

**Theorem 10.** *Let  $\mathcal{R} = \mathcal{O}_{\mathcal{K}_{\mathfrak{g}}} \otimes \mathcal{O}_{\mathcal{K}_{\mathfrak{f}_0}^+} \otimes \dots \otimes \mathcal{O}_{\mathcal{K}_{\mathfrak{f}_{k-1}}^+}$ ,  $\mathfrak{g} = 2^d$  for some  $d \in \mathbb{N}$ , and  $\mathfrak{f}_0, \dots, \mathfrak{f}_{k-1}$  distinct odd primes. Let  $\mathbf{b} = \mathbf{b}_{\mathfrak{g}} \otimes \left( \bigotimes_{i \in [k]} \mathbf{b}_{\mathfrak{f}_i}^+ \right)$ , where  $\mathbf{b}_{\mathfrak{g}}$  is the power basis for  $\mathcal{R}_{\mathfrak{g}}$  and  $\mathbf{b}_{\mathfrak{f}_i}^+$  is a basis for  $\mathcal{R}_{\mathfrak{f}_i}^+$  defined as in Theorem 9. Then,  $\tau(\cdot) = \frac{1}{t} \cdot \text{Trace}(\alpha \cdot (\cdot))$  is inner-product embedding for  $\alpha = \prod_{i \in [k]} \alpha_{\mathfrak{f}_i}$ , where  $t = 2^k \varphi(\mathfrak{g}) \prod_{i \in [k]} \mathfrak{f}_i$ .*

*Proof.* Write  $\mathcal{R}_{\mathfrak{g}}$  for  $\mathcal{O}_{\mathcal{K}_{\mathfrak{g}}}$  and  $\mathcal{R}_i$  for  $\mathcal{O}_{\mathcal{K}_{\mathfrak{f}_i}^+}$  for  $i \in [k]$ . Define  $t_{\mathfrak{f}_i} = 2\mathfrak{f}_i$  and  $t_{\mathfrak{g}} = \varphi(\mathfrak{g})$ . We prove by induction on tensoring consecutive rings  $\mathcal{R}_{\mathfrak{g}}$  and  $\mathcal{R}_{\mathfrak{f}_i} \forall i \in [k]$ , indexed as  $\mathcal{R}_i = \mathcal{O}_{\mathcal{K}_i}$  for  $i \in [h]$ , where  $h = k + 1$  or  $h = k$  (if no power-of-two components). We use  $\mathfrak{h}_i$  for  $i \in [h]$  to iterate over coprime factors of the conductor.

By Remark 7 and Theorem 9  $\mathcal{R}_i$  has an inner-product embedding  $\tau_i$ , i.e.

$$\tau_i((b_i)_m \cdot (\bar{b}_i)_n) = \frac{1}{t_i} \text{Trace}_{\mathcal{K}_i/\mathbb{Q}}(\alpha_i \cdot (b_i)_m \cdot (\bar{b}_i)_n) = \begin{cases} 1 & \text{if } m = n, \\ 0 & \text{if } m \neq n \end{cases} \quad \forall i \in [h].$$

We devise a proof by induction.

First, we define the base case  $(\tilde{\mathfrak{h}}_0, \tilde{t}_0, \tilde{\alpha}_0, \tilde{\mathcal{R}}_0, \tilde{\alpha}_0, \tilde{\mathbf{b}}_0) = (\mathfrak{h}_0, t_0, \alpha_0, \mathcal{R}_0, \alpha_0, \mathbf{b}_0)$ .

Then, for  $i \in [h - 1]$ , define the following inductive steps:

$$\begin{aligned}
\mathfrak{h}_{i+1} &:= \mathfrak{h}_{i+1} \cdot \tilde{\mathfrak{h}}_i & \tilde{t}_{i+1} &:= t_i \cdot \tilde{t}_{i+1} \\
\tilde{\alpha}_{i+1} &:= \alpha_i \cdot \tilde{\alpha}_{i+1} & \tilde{\mathcal{R}}_{i+1} &:= \mathcal{R}_{i+1} \otimes \tilde{\mathcal{R}}_i \\
\tilde{\alpha}_{i+1} &:= \alpha_{i+1} \cdot \tilde{\alpha}_i & \tilde{\mathbf{b}}_{i+1} &:= \mathbf{b}_{i+1} \otimes \tilde{\mathbf{b}}_i
\end{aligned}$$

We want to show that, if  $\tilde{\mathcal{R}}_i$  has an inner-product embedding, then  $\tilde{\mathcal{R}}_{i+1}$  also has an inner-product embedding.

We write

$$\tilde{\mathbf{b}}_i = \{\tilde{b}_{i,0}, \dots, \tilde{b}_{i,\varphi_i-1}\},$$

and

$$\mathbf{b}_{i+1} = \{b_{i+1,0}, \dots, b_{i+1,\varphi_{i+1}-1}\}.$$

Elements of a new basis  $\tilde{\mathbf{b}}_{i+1}$  are uniquely defined as a product of two elements from bases  $\tilde{\mathbf{b}}_i$  and  $\mathbf{b}_{i+1}$ . Consider elements  $b_{i+1,m} \cdot \tilde{b}_{i,r}$  and  $b_{i+1,n} \cdot \tilde{b}_{i,s}$  of a new basis  $\tilde{\mathbf{b}}_{i+1}$ . Due to the coprimality of  $\tilde{\mathfrak{h}}_i$  and  $\mathfrak{h}_{i+1}$ , the tower structure of traces is interchangeable, thus

$$\begin{aligned} & \tilde{\tau}_{i+1} \left( b_{i+1,m} \tilde{b}_{i,r} \cdot \overline{b_{i+1,n} \cdot \tilde{b}_{i,s}} \right) \\ &= \frac{1}{\tilde{t}_{i+1}} \text{Trace}_{\tilde{\mathcal{K}}_{i+1}/\mathbb{Q}} \left( \tilde{\alpha}_{i+1} \cdot b_{i+1,m} \tilde{b}_{i,r} \cdot \overline{b_{i+1,n} \cdot \tilde{b}_{i,s}} \right) \\ &= \frac{1}{\tilde{t}_i} \text{Trace}_{\tilde{\mathcal{K}}_i/\mathbb{Q}} \left( \tilde{\alpha}_i \cdot \tilde{b}_{i,r} \cdot \tilde{b}_{i,s} \right) \cdot \frac{1}{t_{i+1}} \text{Trace}_{\mathcal{K}_{i+1}/\mathbb{Q}} \left( \alpha_{i+1} \cdot b_{i+1,m} \cdot \bar{b}_{i+1,n} \right) \\ &= \tilde{\tau}_i \left( \tilde{b}_{i,r} \cdot \tilde{b}_{i,s} \right) \cdot \tau_{i+1} \left( b_{i+1,m} \cdot \bar{b}_{i+1,n} \right) = \begin{cases} 1 & \text{if } (m, r) = (n, s) \\ 0 & \text{if } (m, r) \neq (n, s) \end{cases} \end{aligned}$$

Finally,

$$(\mathfrak{h}, t, \alpha, \mathcal{R}, \alpha, \mathbf{b}) = \left( \tilde{\mathfrak{h}}_{h-1}, \tilde{t}_{h-1}, \tilde{\alpha}_{h-1}, \tilde{\mathcal{R}}_{h-1}, \tilde{\alpha}_{h-1}, \tilde{\mathbf{b}}_{h-1} \right),$$

which concludes the proof.  $\square$

#### 6.4 Reducing Binariness to Bounded Norm

We show how to reduce the  $\mathbb{Z}$ -relation  $\mathbf{x} \in \{0, 1\}^{m\delta}$  to an  $\mathcal{R}$ -relation natively supported by the succinct arguments presented in Section 5, via the inner-product embedding framework. First, we recall the following elementary fact from [LNP22].

**Proposition 2.** *A vector  $\mathbf{x} \in \mathbb{Z}^{m\delta}$  is binary if and only if  $\langle \mathbf{x}, \mathbf{1}^m - \mathbf{x} \rangle_{\mathbb{Z}} = 0$ .*

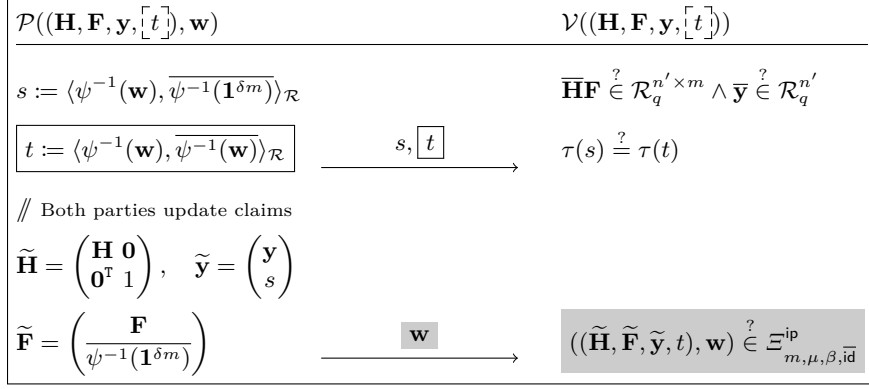
*Proof.* To argue about the “ $\implies$ ” direction this is enough to observe that, for each  $j \in [m\delta]$ , we have  $x_j = 0$  or  $1 - x_j = 0$ . And therefore the sum satisfies  $\sum_{j \in [m\delta]} x_j(1 - x_j) = 0$ . The “ $\impliedby$ ” direction relies on the observation that, for each  $j \in [m\delta]$ ,  $x_j(1 - x_j) \geq 0$  as  $x_j \in \mathbb{Z}$ . Also, if  $x_j \notin \{0, 1\}$ , then  $x_j(1 - x_j) > 0$ . Hence, if for some  $j \in [m\delta]$ ,  $x_j \notin \{0, 1\}$ , then  $\sum_{j \in [m\delta]} x_j(1 - x_j) > 0$ , which is a contradiction.  $\square$

Next, we observe the following equivalence:  $\langle \mathbf{x}, \mathbf{1}^{m\delta} - \mathbf{x} \rangle_{\mathbb{Z}} = 0 \iff \langle \mathbf{x}, \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} - \langle \mathbf{x}, \mathbf{x} \rangle_{\mathbb{Z}} = 0 \iff \langle \mathbf{x}, \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} = \langle \mathbf{x}, \mathbf{x} \rangle_{\mathbb{Z}}$ . This suggests the following reduction:

- (i) The prover sends two claimed values  $s, t \in \mathcal{R}$  supposedly satisfying  $\tau(t) = \tau(s)$
- (ii) The prover then sends a succinct proof for  $\langle \psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{1}^{m\delta})} \rangle_{\mathcal{R}} = s$  and  $\langle \psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{x})} \rangle_{\mathcal{R}} = t$ .

From the identity  $\forall \mathbf{a}, \mathbf{b} \in \mathbb{Z}^{m\delta}$ ,  $\tau(\langle \psi^{-1}(\mathbf{a}), \overline{\psi^{-1}(\mathbf{b})} \rangle_{\mathcal{R}}) = \langle \mathbf{a}, \mathbf{b} \rangle_{\mathbb{Z}}$ , the verifier would be convinced that  $\mathbf{x}$  is indeed binary.

However, there is a subtle issue that, on one hand, the rings  $\mathcal{R}$  considered in this section are of the form displayed in Section 6.3, which are not necessarily equal to  $\mathcal{O}_{\mathcal{K}}$  or  $\mathcal{O}_{\mathcal{K}^+}$  for any cyclotomic field  $\mathcal{K}$ . On the other hand, the succinct arguments constructed in Section 5 are over rings which admit large subtractive sets, for which we only know constructions in  $\mathcal{O}_{\mathcal{K}}$  and  $\mathcal{O}_{\mathcal{K}^+}$ . We therefore need to lift the  $\mathcal{R}$ -relations that we want to prove to some  $\mathcal{O}_{\mathcal{K}}$ -relations (or  $\mathcal{O}_{\mathcal{K}^+}$ -relations, but we focus on the former) with  $\mathcal{R} \subseteq \mathcal{O}_{\mathcal{K}}$ , while ensuring that the prover cannot cheat by using a witness over  $\mathcal{O}_{\mathcal{K}}$ . To do this, we need a lemma which allows viewing  $\mathcal{O}_{\mathcal{K}}$  as an  $\mathcal{R}$ -module in such a way that the geometry of  $\mathcal{O}_{\mathcal{K}}$  is respected. A precise statement is given in Lemma 8.



**Fig. 9.** Protocol  $\boxed{\Pi_{\tau}^{\text{lin-bin}}}$  or  $\boxed{\Pi_{\tau}^{\text{ip-bin}}}$ , a reduction from  $\boxed{\Xi_{m, \mu, \beta}^{\text{lin}} \cap \Xi_m^{\text{bin}}}$  or  $\boxed{\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}} \cap \Xi_m^{\text{bin}}}$  to  $\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}}$ . The marked parts are only sent / checked when the protocol is used as a proof of knowledge. As a reduction of knowledge, they are omitted.

We next formally define the binariness relation which ignores the statement and simply checks that the witness is a binary vector.

$$\Xi_m^{\text{bin}} := \{(\text{stmt}, \mathbf{w}) : \text{stmt} \in \{0, 1\}^*; \mathbf{w} \in \mathcal{R}^m; \psi(\mathbf{w}) \in \{0, 1\}^{m\delta} \}.$$

In Fig. 9, we present two similar reductions of knowledge  $\Pi^{\text{lin-bin}}$  and  $\Pi^{\text{ip-bin}}$  from  $\Xi^{\text{lin}} \cap \Xi^{\text{bin}}$  or  $\Xi^{\text{ip}} \cap \Xi^{\text{bin}}$  to  $\Xi^{\text{ip}}$ , respectively. Note that, when reducing  $\Xi^{\text{ip}} \cap \Xi^{\text{bin}}$  to  $\Xi^{\text{ip}}$ , the inner product  $t = \langle \psi^{-1}(\mathbf{x}), \overline{\psi^{-1}(\mathbf{x})} \rangle_{\mathcal{R}}$  is already included as part of the statement, and thus the prover does not need to send it. The formal result is stated in Theorem 11, whose proof relies on Lemma 8 stated immediately after.

**Theorem 11.** *Let  $\mathcal{R} = \mathcal{O}_{\mathcal{K}_g} \otimes \mathcal{O}_{\mathcal{K}_{p_0}^+} \otimes \dots \otimes \mathcal{O}_{\mathcal{K}_{p_{k-1}}^+}$ ,  $\mathfrak{g} = 2^d$  for some  $d \in \mathbb{N}$ , and  $f_0, \dots, f_{k-1}$  distinct odd primes. Let  $\tau$  be an inner-product embedding over  $\mathcal{R}$ . The protocol  $\Pi_{\tau}^{\text{lin-bin}}$  (resp.  $\Pi_{\tau}^{\text{ip-bin}}$ ) is a perfectly correct reduction of knowledge from  $\Xi_{m, \mu, \beta}^{\text{lin}} \cap \Xi_m^{\text{bin}}$  (resp.  $\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}} \cap \Xi_m^{\text{bin}}$ ) over  $\mathcal{R}$  to  $\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}}$  over  $\mathcal{O}_{\mathcal{K}}$ . There exists a constant  $c_f$  such that it is relaxed knowledge sound from  $\Xi_{m, \mu, \beta}^{\text{lin}} \vee \text{sis} \cap \Xi_m^{\text{bin}}$  (resp.  $\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}} \vee \text{sis} \cap \Xi_m^{\text{bin}}$ ) over  $\mathcal{R}$  to  $\Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}}$  over  $\mathcal{O}_{\mathcal{K}}$  if  $2^k c_f \cdot \varphi^{5/2} \cdot \beta \leq \beta^{\text{sis}}$ .*

*Proof.* For perfect completeness, consider  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}) \in \Xi_{m, \mu, \beta}^{\text{lin}} \cap \Xi_m^{\text{bin}}$  over  $\mathcal{R}$ . We have  $\psi(\mathbf{w}) \in \{0, 1\}^{m\delta}$ . By Proposition 2 and the discussion immediately after, it holds that  $\langle \mathbf{w}, \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} = \langle \mathbf{w}, \mathbf{w} \rangle_{\mathbb{Z}}$ . Since  $\tau$  is an inner-product embedding over  $\mathcal{R}$ , it holds that

$$\begin{aligned} \tau(s) &= \tau(\langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{1}^{\delta m})} \rangle_{\mathcal{R}}) = \langle \mathbf{w}, \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}}, \text{ and} \\ \tau(t) &= \tau(\langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{w})} \rangle_{\mathcal{R}}) = \langle \mathbf{w}, \mathbf{w} \rangle_{\mathbb{Z}}, \end{aligned}$$

and thus  $\tau(s) = \tau(t)$ . Furthermore, since  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}) \in \Xi_{m, \mu, \beta}^{\text{lin}}$ , we have  $s = \langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{1}^{\delta m})} \rangle_{\mathcal{R}}$ , and  $t = \langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{w})} \rangle_{\mathcal{R}}$ . Therefore,  $((\tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}, t), \mathbf{w}) \stackrel{?}{\in} \Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}}$  over  $\mathcal{R}$ . Since  $\mathcal{R} \subseteq \mathcal{O}_{\mathcal{K}}$ , the claim follows.

For perfect relaxed knowledge soundness, suppose that  $((\tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}, t), \mathbf{w}) \stackrel{?}{\in} \Xi_{m, \mu, \beta, \bar{\text{id}}}^{\text{ip}}$  over  $\mathcal{O}_{\mathcal{K}}$ . If  $\mathbf{w} \in \mathcal{R}^m$ , then we have  $s = \langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{1}^{\delta m})} \rangle_{\mathcal{R}}$  and  $t = \langle \psi^{-1}(\mathbf{w}), \overline{\psi^{-1}(\mathbf{w})} \rangle_{\mathcal{R}}$ . Since  $\tau(s) = \tau(t)$ , reversing the above argument gives  $\psi(\mathbf{w}) \in \{0, 1\}^{m\delta}$ . Thus  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}) \in \Xi_{m, \mu, \beta}^{\text{lin}} \cap \Xi_m^{\text{bin}}$  over  $\mathcal{R}$  as desired.

If  $\mathbf{w} \in \mathcal{O}_{\mathcal{K}}^m \setminus \mathcal{R}^m$ , then we can express  $\mathbf{w}$  as a linear combination of  $\mathcal{R}$ -vectors. Let  $\hat{\mathbf{w}}$  be the coefficient of any basis element other than 1. Note that  $\overline{\mathbf{H}}\mathbf{F}\hat{\mathbf{w}} = \mathbf{0}^{n'}$  mod  $q$ . Moreover, by Lemma 8, we have  $\|\sigma(\hat{\mathbf{w}})\|_2 \leq c_f \cdot \sqrt{m} \cdot \varphi^{5/2} \cdot \beta \leq \beta^{\text{sis}}$ , i.e. a vSIS break.

The argument for  $\Pi_{\tau}^{\text{ip-bin}}$  is almost verbatim, except that  $t$  is given as part of the statement rather than being sent by the prover.  $\square$

**Lemma 8.** *If  $f$  be an odd prime, then  $\mathcal{O}_{\mathcal{K}}$  can be seen as an  $\mathcal{O}_{\mathcal{K}^+}$ -module with the basis  $\{1, \zeta\}$ . More generally, let  $\mathcal{R} = \mathcal{O}_{\mathcal{K}_{\mathfrak{g}}} \otimes \mathcal{O}_{\mathcal{K}_{f_0}^+} \otimes \dots \otimes \mathcal{O}_{\mathcal{K}_{f_{k-1}}^+}$  where  $\mathfrak{g} = 2^d$  for some  $d \in \mathbb{N}$  and  $f_0, \dots, f_{k-1}$  are distinct odd primes. Let  $\mathfrak{f} := \mathfrak{g} \prod_{i \in [k]} f_i$ . Then  $\mathcal{O}_{\mathcal{K}_{\mathfrak{f}}}$  is an  $\mathcal{R}$ -module with the basis  $\bigotimes_{i \in [k]} (1, \zeta_{f_i})$ .*

*Furthermore, let  $\mathbf{x} \in \mathcal{O}_{\mathcal{K}_{\mathfrak{f}}}^m$  be expressed as a  $\mathcal{R}^m$ -combination of  $\bigotimes_{i \in [k]} (1, \zeta_{f_i})$ . There exists a constant  $c_{\mathfrak{f}}$  (dependent only on  $\mathfrak{f}$ ), such that if  $\|\sigma(\mathbf{x})\|_2 \leq \beta$ , then each  $\mathcal{R}$ -coefficient  $\hat{\mathbf{x}}$  of  $\mathbf{x}$  satisfies  $\|\sigma(\hat{\mathbf{x}})\|_2 \leq \beta'$ , where  $\beta' \leq 2^k c_{\mathfrak{f}} \cdot \varphi^{5/2} \cdot \beta$ .*

*Proof. First part.* We first consider the simple case where  $f$  is an odd prime. Consider the  $\mathbb{Z}$ -basis  $\mathbf{b}^+ = \{b_0^+, \dots, b_{\varphi/2-1}^+\} = \{1, \zeta + \zeta^{-1}, \dots, \zeta^{\varphi/2-1} + \zeta^{1-\varphi/2}\}$  of  $\mathcal{O}_{\mathcal{K}^+}$ . We show that  $(1, \zeta) \otimes \mathbf{b}^+$  is a  $\mathbb{Z}$ -basis of  $\mathcal{O}_{\mathcal{K}}$ , implying that  $\mathcal{O}_{\mathcal{K}}$  is an  $\mathcal{O}_{\mathcal{K}^+}$ -module with the basis  $\{1, \zeta\}$ . Consider the ‘‘balanced power basis’’

$$\begin{aligned} \mathbf{b} &= \{b_{-\varphi/2+1}, \dots, b_{-1}, b_0, \dots, b_{\varphi/2}\} \\ &= \{\zeta^{-\varphi/2+1}, \dots, \zeta^{-1}, 1, \dots, \zeta^{\varphi/2}\}. \end{aligned}$$

We prove by induction that  $b_{-i}$  and  $b_{i+1}$  can be expressed as a  $\{-1, 0, 1\}$ -combination of elements from  $(1, \zeta) \otimes \mathbf{b}^+$  for all  $i \in [\varphi/2]$ .

For  $i = 0$ , we observe that  $b_0 = b_0^+$  and  $b_1 = b_0^+ \cdot \zeta$ . Now, suppose the induction hypothesis holds for  $i \leq k$  for some  $\ell \in [\varphi/2]$ , i.e.  $b_{-i}$  and  $b_{i+1}$  are constructed for all  $i \in [k]$ . Our goal is to obtain  $b_{-k}$  and  $b_{k+1}$ . Observe that

$$\begin{aligned} b_k^+ &= b_k + b_{-k}, \\ b_k^+ \cdot \zeta &= (b_k + b_{-k}) \cdot \zeta = b_{k+1} + b_{-k+1}, \\ b_{-k} &= b_k^+ - b_k, \\ b_{k+1} &= b_k^+ \cdot \zeta - b_{-k+1}. \end{aligned}$$

The claim then follows from the induction hypothesis. The above directly generalises for the tensor rings by arguing about each factor ring independently.

*‘‘Furthermore’’ part.* For simplicity, we first consider the case where  $f$  is prime. From the previous part of the proof, we know that there exists a  $\mathbb{Z}$ -basis  $\mathbf{b}$  of  $\mathcal{O}_{\mathcal{K}}$  which can be expressed as a  $\{-1, 0, 1\}$ -combination of elements in  $(1, \zeta) \otimes \mathbf{b}^+$ . We can write

$$b_i = \sum_{j \in [\varphi/2]} s_{i,j} b_j^+ + \sum_{j \in [\varphi/2]} t_{i,j} b_j^+ \cdot \zeta,$$

where  $s_{i,j}, t_{i,j} \in \{-1, 0, 1\}$ .

Consider  $\mathbf{x} = (x_0, \dots, x_{m-1})$  for any  $\mathbf{x} \in \mathcal{O}_{\mathcal{K}_{\mathfrak{f}}}^m$ . Then, we can write

$$x_i = \sum_{j \in [\varphi]} \tilde{x}_{i,j} b_j = \sum_{j \in [\varphi], \ell \in [\varphi/2]} \tilde{x}_{i,j} (s_{j,\ell} + t_{j,\ell} \zeta) b_{\ell}^+,$$

where  $\tilde{x}_{i,j} \in \mathbb{Z}^m$ . Consider  $\|\sigma(x_i)\|_2 = \beta_i$ , by applying the same argument as in [DPSZ11, Theorem 7] but with the basis  $\mathbf{b}$ , we conclude that there exists a constant  $c_{\mathfrak{f}}$  (dependent only on  $\mathfrak{f}$ ) such that  $\|\psi(x_i)\|_{\infty} \leq c_{\mathfrak{f}} \beta_i$ . Corollary 2 discusses the constant  $c_{\mathfrak{f}}$ .

Let  $\hat{x}_{i,\ell} = \sum_{j \in [\varphi]} \tilde{x}_{i,j} (s_{j,\ell} + t_{j,\ell} \zeta)$ . Then,

$$x_i = \sum_{j \in [\varphi/2]} \hat{x}_{i,j} b_j^+.$$

We observe that as  $\|\psi(\tilde{x}_{i,j})\|_{\infty} \leq c_{\mathfrak{f}} \beta_i$ , then  $\|\psi(\hat{x}_{i,\ell})\|_{\infty} \leq 2\varphi c_{\mathfrak{f}} \beta_i$  and  $\|\sigma(\hat{x}_{i,\ell})\|_2 \leq 2\varphi^{5/2} c_{\mathfrak{f}} \beta_i$  due to the norm conversion.

Eventually, consider the norm of  $\tilde{\mathbf{x}}_j$ , i.e.

$$\|\sigma(\tilde{\mathbf{x}}_j)\|_2 \leq \sqrt{\sum_{i \in [m]} (\beta_i 2\varphi^{5/2} c_{\mathfrak{f}})^2} = 2\varphi^{5/2} c_{\mathfrak{f}} \sqrt{\sum_{i \in [m]} \beta_i^2} = 2\varphi^{5/2} c_{\mathfrak{f}} \|\sigma(\mathbf{x})\|_2 \leq 2\varphi^{5/2} c_{\mathfrak{f}} \beta.$$

To argue about the composite case, we consider bases  $\mathbf{b} = \mathbf{b}^{(1)} \otimes \mathbf{b}^{(2)}$ , where  $\mathbf{b}^{(1)}$  and  $\mathbf{b}^{(2)}$  are bases of the cyclotomic rings with prime conductors. Let  $b_{i,j} = b_i \cdot b_j$  and  $b_{i,j}^+ = b_i^+ \cdot b_j^+$ . Let  $\varphi = \varphi^{(1)} \cdot \varphi^{(2)}$  and  $\zeta = \zeta^{(1)} \cdot \zeta^{(2)}$  defined analogously. Then, we write:

$$\begin{aligned} x_i &= \sum_{j^{(1)} \in [\varphi^{(1)}], j^{(2)} \in [\varphi^{(2)}]} \tilde{x}_{i,j^{(1)},j^{(2)}} b_{\ell^{(1)}}^{(1)} b_{\ell^{(2)}}^{(2)} \\ &= \sum_{\substack{j^{(1)} \in [\varphi^{(1)}] \\ j^{(2)} \in [\varphi^{(2)}] \\ \ell^{(1)} \in [\varphi^{(1)}/2] \\ \ell^{(2)} \in [\varphi^{(2)}/2]}} \tilde{x}_{i,j^{(1)},j^{(2)}} b_{\ell^{(1)}}^{(1)} b_{\ell^{(2)}}^{(2)} = \sum_{\substack{j^{(1)} \in [\varphi^{(1)}] \\ j^{(2)} \in [\varphi^{(2)}] \\ \ell^{(1)} \in [\varphi^{(1)}/2] \\ \ell^{(2)} \in [\varphi^{(2)}/2]}} \tilde{x}_{i,j^{(1)},j^{(2)}} (s_{j^{(1)},\ell^{(1)}}^{(1)} + t_{j^{(1)},\ell^{(1)}}^{(1)} \cdot \zeta^{(1)}) (s_{j^{(2)},\ell^{(2)}}^{(2)} + t_{j^{(2)},\ell^{(2)}}^{(2)} \cdot \zeta^{(2)}) \cdot b_{\ell^{(1)},\ell^{(2)}}^+ \end{aligned}$$

Let  $\hat{x}_{i,\ell^{(1)},\ell^{(2)}} = \sum_{j^{(1)} \in [\varphi^{(1)}], j^{(2)} \in [\varphi^{(2)}]} \tilde{x}_{i,j^{(1)},j^{(2)}} (s_{j^{(1)},\ell^{(1)}}^{(1)} + t_{j^{(1)},\ell^{(1)}}^{(1)} \cdot \zeta^{(1)}) (s_{j^{(2)},\ell^{(2)}}^{(2)} + t_{j^{(2)},\ell^{(2)}}^{(2)} \cdot \zeta^{(2)})$ . Then, We observe that as  $\|\psi(\tilde{x}_{i,j^{(1)},j^{(2)}})\|_\infty \leq c_f \beta_i$ , then  $\|\psi(\hat{x}_{i,\ell^{(1)},\ell^{(2)}})\|_\infty \leq 4\varphi c_f \beta_i$  and continue the reasoning as in the base case. Clearly, the argument extends for terson rings of more than two prime rings.  $\square$

## 7 Packed $\mathbb{Z}$ -Inner Products via CRT Embedding

The idea of embedding  $\mathbb{Z}$ -relations into  $\mathcal{R}$ -relations via the CRT embedding is well-established (e.g. [BS23, LNS20]). However, an obstacle to applying this to lattice-based succinct arguments is the lack of a succinct-verifier argument for proving the consistency of two vectors related via the coefficient and the CRT embeddings.

In this section, we first recall the method of embedding  $\mathbb{Z}$ -relations into  $\mathcal{R}$ -relations via the CRT embedding. Then, by exploiting the fine-grained tower structure of cyclotomic rings with smooth conductors, we provide a verifier-succinct argument for proving the consistency between the coefficient and the CRT embeddings. Throughout this section, we assume that  $\mathcal{R} = \mathbb{Z}[\zeta_f]$  is a cyclotomic ring of degree  $\varphi$ , and  $p \in \mathbb{N}$  is a rational prime which splits completely over  $\mathcal{R}$ .

### 7.1 Embedding $\mathbb{Z}_p$ -inner products into $\mathcal{R}_p$ -inner products

To begin, let us write  $\text{CRT}_p : \mathcal{R} \rightarrow \mathbb{Z}^\varphi$ , for the invertible  $\mathbb{Z}$ -linear transform which maps a ring element  $x \in \mathcal{R}$  to its Chinese remainder representation modulo each prime ideal dividing  $p$ . Note that we are viewing  $\text{CRT}_p$  as a  $\mathbb{Z}$ -linear map rather than a  $\mathbb{Z}_p$ -linear map, and we will write  $\text{mod } p$  explicitly when reducing modulo  $p$ . We extend the notation naturally to vectors, i.e. for  $\mathbf{x} = (x_i)_{i \in [m]} \in \mathcal{R}^m$  we define  $\text{CRT}_p(\mathbf{x}) = (\text{CRT}_p(x_i))_{i \in [m]}$ .

It is well-known that addition and multiplication in the CRT domain is component-wise. Using this property, there exists a natural method of embedding  $\mathbb{Z}_p$ -inner products into  $\mathcal{R}_p$ -inner-products, as summarised in Proposition 3. The proof is trivial and thus omitted.

**Proposition 3.** *Let  $\mathcal{R} = \mathbb{Z}[\zeta_f]$  be a cyclotomic ring of degree  $\varphi$  and  $p \in \mathbb{N}$  be a rational prime which fully splits over  $\mathcal{R}$ . Let  $\tau_p : \mathcal{R} \rightarrow \mathbb{Z}$  be defined as  $\tau_p(z) := \langle \mathbf{1}^\varphi, \text{CRT}_p(z) \rangle$ . For any  $\mathbf{x} = (x_i)_{i \in [m]}, (\mathbf{y}_i)_{i \in [m]} \in \mathbb{Z}^{m\varphi}$ , it holds that*

$$\tau_p(\langle \text{CRT}_p^{-1}(\mathbf{x}), \text{CRT}_p^{-1}(\mathbf{y}) \rangle_{\mathcal{R}}) = \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} \text{ mod } p.$$

Using Proposition 3, a prover is already able to succinctly prove that certain  $\mathbb{Z}_p$ -inner product relations hold using the succinct arguments provided in Section 5, provided that the application allows the witness vectors to be committed in their  $\text{CRT}_p^{-1}(\cdot)$  form. A bit more concretely, consider a toy example where the prover wishes to prove that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = z \text{ mod } p$  for some public value  $z \in \mathbb{Z}_p$ , where  $p$  is sufficiently shorter than the modulus  $q$  used in the argument system. It performs the following procedures:

- Compute  $\hat{z} := \langle \text{CRT}_p^{-1}(\mathbf{x}), \text{CRT}_p^{-1}(\mathbf{y}) \rangle_{\mathcal{R}} \text{ mod } p$  and send it to the verifier.
- Let  $\hat{\mathbf{x}} := \text{CRT}_p^{-1}(\mathbf{x})$  and  $\hat{\mathbf{y}} := \text{CRT}_p^{-1}(\mathbf{y}) \text{ mod } p$ .
- Find  $\mathbf{r} \in \mathcal{R}^m$  such that  $\hat{z} = \langle \hat{\mathbf{x}}, \hat{\mathbf{y}} \rangle_{\mathcal{R}} + p \cdot \mathbf{r}$ .
- Commit to  $(\hat{\mathbf{x}}, \hat{\mathbf{y}}, \mathbf{r})$ .



- Provide a proof that  $(\widehat{\mathbf{x}}, \widehat{\mathbf{y}}, \mathbf{r})$  satisfies  $\widehat{z} = \langle \widehat{\mathbf{x}}, \widehat{\mathbf{y}} \rangle_{\mathcal{R}} + p \cdot \mathbf{r}$ .

In turn, the verifier checks that  $\tau_p(\widehat{z}) = z \bmod p$  and the proof for  $\widehat{z} = \langle \widehat{\mathbf{x}}, \widehat{\mathbf{y}} \rangle_{\mathcal{R}} + p \cdot \mathbf{r}$  is valid. If both checks go through, then by the soundness of the argument system the verifier would be convinced that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = z \bmod p$  (for some  $\mathbf{x}, \mathbf{y}$  which satisfy  $(\mathbf{x}, \mathbf{y}) = (\text{CRT}_p(\widehat{\mathbf{x}}), \text{CRT}_p(\widehat{\mathbf{y}})) \bmod p$ ).

## 7.2 Lifting to $\mathbb{Z}$ and $\mathcal{R}$

In case the prover wishes to prove that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = z$  without reduction modulo  $p$ , and/or the application postulates that the witness vectors are committed in  $\psi^{-1}$  form, then the prover needs to additionally perform the following:

- Write  $\widetilde{\mathbf{x}} := \psi^{-1}(\mathbf{x})$  and  $\widetilde{\mathbf{y}} := \psi^{-1}(\mathbf{y})$ .
- Find  $\mathbf{r}, \widetilde{\mathbf{s}} \in \mathcal{R}^m$  such that

$$\widetilde{\mathbf{x}} = \psi^{-1}(\text{CRT}_p(\widehat{\mathbf{x}})) + p \cdot \mathbf{r}, \quad (7)$$

$$\widetilde{\mathbf{y}} = \psi^{-1}(\text{CRT}_p(\widehat{\mathbf{y}})) + p \cdot \widetilde{\mathbf{s}}. \quad (8)$$

- Further commit to  $(\widetilde{\mathbf{x}}, \widetilde{\mathbf{y}}, \mathbf{r}, \widetilde{\mathbf{s}})$ .
- Provide a proof that  $\|\sigma(\widetilde{\mathbf{x}})\|_2$  and  $\|\sigma(\widetilde{\mathbf{y}})\|_2$  are short.
- Provide a proof that Eqs. (7) and (8) hold.

From Eqs. (7) and (8), the verifier would be convinced that

$$\psi(\widetilde{\mathbf{x}}) = \text{CRT}_p(\widehat{\mathbf{x}}) \bmod p \quad \text{and} \quad \psi(\widetilde{\mathbf{y}}) = \text{CRT}_p(\widehat{\mathbf{y}}) \bmod p.$$

Combined with the previous guarantee that  $\langle \text{CRT}_p(\widehat{\mathbf{x}}), \text{CRT}_p(\widehat{\mathbf{y}}) \rangle_{\mathbb{Z}} = z \bmod p$ , the verifier would be convinced that

$$\langle \psi(\widetilde{\mathbf{x}}), \psi(\widetilde{\mathbf{y}}) \rangle_{\mathbb{Z}} = z \bmod p.$$

With the proof of  $\|\sigma(\widetilde{\mathbf{x}})\|_2$  and  $\|\sigma(\widetilde{\mathbf{y}})\|_2$  being short, it must be the case that  $\|\psi(\widetilde{\mathbf{x}})\|_{\infty}$  and  $\|\psi(\widetilde{\mathbf{y}})\|_{\infty}$  are also short. Provided that these norms are small enough relative to  $p$ , the reduction modulo  $p$  has no effect, and thus we arrive at

$$\langle \psi(\widetilde{\mathbf{x}}), \psi(\widetilde{\mathbf{y}}) \rangle_{\mathbb{Z}} = z.$$

Next, we discuss how to succinctly instantiate the above protocol, in particular the arguments for Eqs. (7) and (8), using tools developed in Section 5.

## 7.3 Computing CRT via Automorphisms

In the above, we established that the method of embedding  $\mathbb{Z}$ -inner products via CRT requires proving consistency between the  $\text{CRT}_p^{-1}(\cdot)$  and  $\psi^{-1}(\cdot)$  encodings of the witness vectors. Specifically, we would like to design a succinct argument for arguing that

$$\widetilde{\mathbf{x}} = \psi^{-1}(\text{CRT}_p(\widehat{\mathbf{x}})) \bmod p$$

where  $\widetilde{\mathbf{x}}, \widehat{\mathbf{x}} \in \mathcal{R}^m$  are committed vectors. We note that existing protocols for proving such correspondence treat  $\psi^{-1} \circ \text{CRT}_p$  as a generic  $\mathbb{Z}$ -linear map and prove the correspondence as an unstructured system of linear equations over  $\mathbb{Z}$ . Instead, we would like to exploit the tensor structure of the  $\psi^{-1} \circ \text{CRT}_p$  map for carefully chosen rings  $\mathcal{R}$ , and the fact that any  $\mathbb{Z}$ -linear map can be expressed as a linear combination of automorphisms in  $\text{Gal}(\mathcal{K}/\mathbb{Q})$  with  $\mathcal{R}$  coefficients.

Motivated by the above, the goal of this subsection is to prove Theorem 12, which states that, for  $\mathcal{R}$  with a smooth conductor, the  $\psi^{-1} \circ \text{CRT}_p$  map can be expressed as the composition of a few succinct linear combinations of automorphisms in  $\text{Gal}(\mathcal{K}/\mathbb{Q})$  with  $\mathcal{R}$  coefficients. To prove this theorem, we will make use of two elementary lemmas. In Lemma 11, we prove an elementary fact that, if  $L/K$  is a Galois extension, then any  $K$ -linear map  $f : L \rightarrow L$  can be expressed as an  $L$ -linear combination of  $\text{Gal}(L/K)$ . Then, in Lemma 12, we prove an analogous lemma for  $\mathcal{O}_K/p\mathcal{O}_K$ -linear map  $f : \mathcal{O}_L/p\mathcal{O}_L \rightarrow \mathcal{O}_L/p\mathcal{O}_L$ , if  $L$  is cyclotomic and has conductor less than  $p$ , where  $p$  is a rational prime. Using these two results, we arrive at the following theorem.

**Theorem 12.** Let  $\mathcal{R}$  be a cyclotomic ring<sup>18</sup> with a  $w$ -smooth conductor  $\mathfrak{f}$ . Then, the transformation  $\psi^{-1} \circ \text{CRT}_p : \mathcal{R} \rightarrow \mathcal{R}$  can be expressed as the succinct composition of  $\mathcal{R}$ -linear combinations of at most  $w$  automorphisms over  $\mathcal{R}$ . Formally, there exists  $t \in O(\log \mathfrak{f})$ ,  $h_i \leq w$ ,  $s_{i,j} \in \mathcal{R}$ , and  $\alpha_{i,j} \in \text{Gal}(\mathcal{K}/\mathbb{Q})$  for all  $i \in [t]$  and  $j \in [h_i]$ , such that

$$(\psi^{-1} \circ \text{CRT}_p)(\cdot) = \bigcirc_{i \in [t]} \sum_{j \in [h_i]} s_{i,j} \alpha_{i,j}(\cdot) \pmod{p},$$

where  $\bigcirc$  denotes function composition.

*Proof.* Let  $\mathcal{K}$  denote the  $\mathfrak{f}$ -th cyclotomic field. Since  $\mathfrak{f}$  is  $w$ -smooth,  $\mathcal{K}/\mathbb{Q}$  can be decomposed into a tower of  $t \leq O(\log \mathfrak{f})$  Galois extensions where each step of the extension is of degree at most  $h_i \leq w$ . Let  $L/K$  denote the  $i$ -th step of the tower of extensions. Correspondingly, the map  $(\psi^{-1} \circ \text{CRT}_p \pmod{p})$  can be decomposed as a composition

$$\psi^{-1} \circ \text{CRT}_p = \bigcirc_{i \in [t]} \widehat{f}_i$$

where  $\widehat{f}_i$  is obtained by lifting an  $\mathcal{O}_K/p\mathcal{O}_K$ -linear map  $f_i : \mathcal{O}_L/p\mathcal{O}_L \rightarrow \mathcal{O}_L/p\mathcal{O}_L$  to  $\mathcal{O}_K/p\mathcal{O}_K$ . By Lemma 12,  $f_i$  can be expressed as an  $\mathcal{O}_L/p\mathcal{O}_L$ -linear combination of  $\text{Gal}(L/K)$  which contains at most  $h_i \leq w$  elements. Correspondingly,  $\widehat{f}_i$  can be expressed as an  $\mathcal{O}_K/p\mathcal{O}_K$ -linear combination of  $\text{Gal}(\mathcal{K}/\mathbb{Q})$  which contains at most  $h_i \leq w$  elements. The theorem thus follows.  $\square$

#### 7.4 $\Pi^{\text{eip}+}$ : Extended Inner-product relation

The  $\Xi^{\text{ip}}$  relation defined in Section 5.7 asserts a single constraint on the self-inner-product of the entire witness. In preparation for our CRT-based embedding for  $\mathbb{Z}$ -inner-products to be presented in Section 7, we need a slightly extended relation which captures multiple inner-product relations between different blocks of the witness vector, which is now interpreted as a block vector. Formally, we define the ‘‘extended inner-product’’ relation  $\Xi^{\text{eip}}$  below.

$$\Xi_{m, n^{\text{out}}, \beta, \beta^{\text{sis}}}^{\text{eip} \vee \text{sis}} := \left\{ \left( (\iota_{\text{ip}}, \iota_{\text{ip-in}}, \mathbf{H}, \mathbf{F}, \mathbf{y}, \mathbf{t}), \mathbf{w} = (\mathbf{w}_k)_{k \in [n^{\text{blk}}]} \right) : \right. \\ \left. \begin{array}{l} \mathbf{H} \in \mathcal{R}_q^{n^{\text{out}} \times n}; \mathbf{F} \in \mathcal{R}_q^{n \times (d^{\otimes \mu} \cdot n^{\text{blk}})} \subseteq \mathcal{R}_q^{n \times m}; \mathbf{y} \in \mathcal{R}_q^{n^{\text{out}}}; \mathbf{t} \in \mathcal{R}_q^{n^{\text{ip}}} \\ \left\{ \begin{array}{l} \iota_{\text{ip}} \in \{1, -1\}^{[n^{\text{ip}}]}; \iota_{\text{ip-in}} \in [n^{\text{blk}}]^{[n^{\text{ip}}]} \\ \|\mathbf{w}\| \leq \beta; \mathbf{H}\mathbf{F}\mathbf{w} = \mathbf{H}\mathbf{y} \pmod{q} \\ \langle \mathbf{w}_{\iota_{\text{ip-in}}(k)}, \alpha_{\iota_{\text{ip}}(k)}(\mathbf{w}_{\iota_{\text{ip-in}}(k)}) \rangle = t_k \pmod{q} \quad \forall k \in [n^{\text{ip}}] \end{array} \right\} \end{array} \right\} \quad \text{or} \quad \left\{ \begin{array}{l} \|\mathbf{w}\| \leq \beta^{\text{sis}} \\ \overline{\mathbf{H}\mathbf{F}\mathbf{w}} = \mathbf{0}_{\overline{n}} \pmod{q} \end{array} \right\}.$$

Note that  $\Xi^{\text{eip}}$  implicitly has two additional parameters: number of blocks  $n^{\text{blk}}$  and number of inner-product relations  $n^{\text{ip}}$ . Compared to  $\Xi^{\text{ip}}$ , a statement in  $\Xi^{\text{eip}}$  contains additionally two index maps  $\iota_{\text{ip}} : [n^{\text{ip}}] \rightarrow \{1, -1\}$  and  $\iota_{\text{ip-in}} : [n^{\text{ip}}] \rightarrow [n^{\text{blk}}]$ , and the inner product image  $t$  is replaced by a vector  $\mathbf{t} \in \mathcal{R}_q^{n^{\text{ip}}}$ . Furthermore, the single inner-product relation in  $\Xi^{\text{ip}}$  is replaced with

$$\langle \mathbf{w}_{\iota_{\text{ip-in}}(k)}, \alpha_{\iota_{\text{ip}}(k)}(\mathbf{w}_{\iota_{\text{ip-in}}(k)}) \rangle = t_k \pmod{q} \quad \forall k \in [n^{\text{ip}}].$$

Recall that  $\Pi^{\text{ip}+}$  is a reduction of knowledge from  $\Xi^{\text{ip}}$  to  $\Xi^{\text{lin}}$ . Similarly, we can also construct a reduction of knowledge  $\Pi^{\text{eip}+} : \Xi^{\text{eip}} \rightarrow \Xi^{\text{lin}}$ , by modifying the  $\Pi^{\text{ip}}$  (Fig. 8) part of  $\Pi^{\text{ip}+}$  into  $\Pi^{\text{eip}}$ .

Since  $\Pi^{\text{ip}}$  and  $\Pi^{\text{ip}+}$  are already quite notation heavy, and their generalisations to  $\Pi^{\text{eip}}$  and  $\Pi^{\text{eip}+}$  are straightforward, we omit a formal description but instead highlight the differences between  $\Pi^{\text{eip}}$  and  $\Pi^{\text{ip}}$  below.

**The main protocol** The main protocol  $\Pi^{\text{eip}}$  differs from  $\Pi^{\text{ip}}$  in the following:

- (i) The protocol runs in parallel for  $j \in [n^{\text{ip}}]$ :
  - (i) The witness  $\mathbf{w}$  used for obtaining  $\mathbf{v}$  (now, denoted as  $\widehat{\mathbf{v}}_j$ ) is replaced by  $\widehat{\mathbf{w}}_j := \mathbf{e}_{\iota_{\text{ip-in}}(j)} \otimes \mathbf{w}_{\iota_{\text{ip-in}}(j)}$ .

<sup>18</sup>The technique should apply for the ring of integers of any Abelian number field with a known  $w$ -smooth tower structure, but a formal proof is out of scope.

(ii) Similarly, matrix  $\mathbf{E}$  is replaced by

$$\hat{\mathbf{E}}_j := \mathbf{e}_{\iota_{\text{ip-in}}(j)}^T \otimes \begin{pmatrix} 1 & \xi & \dots & \xi^{m/n^{\text{blk}}-1} \\ 1 & \xi^{-1} & \dots & \xi^{-(m/n^{\text{blk}}-1)} \\ 1 & 0 & \dots & 0 \end{pmatrix}.$$

(ii) The new claims both parties compute are:

$$\begin{aligned} \tilde{\mathbf{H}} &:= \begin{pmatrix} \mathbf{H} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n^{\text{ip}} \cdot 3} \end{pmatrix}, & \tilde{\mathbf{F}} &:= \begin{pmatrix} \mathbf{F} \\ \sum_{j \in [n^{\text{ip}}]} \mathbf{e}_j \otimes \hat{\mathbf{E}}_j \end{pmatrix}, \\ \tilde{\mathbf{y}} &:= \begin{pmatrix} \mathbf{y} \\ \sum_{j \in [n^{\text{ip}}]} \mathbf{e}_j \otimes \mathbf{y}_{\hat{\mathbf{E}}_j} \end{pmatrix}, & \tilde{\mathbf{y}}'_i &:= \begin{pmatrix} \mathbf{y}'_i \\ \sum_{j \in [n^{\text{ip}}]} \mathbf{e}_j \otimes \mathbf{y}'_{\hat{\mathbf{E}}_j, i} \end{pmatrix}, \end{aligned}$$

where  $\mathbf{y}_{\hat{\mathbf{E}}_j} := \hat{\mathbf{E}}_j \mathbf{w}$  and  $\mathbf{y}'_{\hat{\mathbf{E}}_j, i} := \hat{\mathbf{E}}_j \hat{\mathbf{v}}_{j, i}$ .

(iii) Further, the verifications are replaced by

$$\begin{aligned} & ((\tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}), \mathbf{w}) \stackrel{?}{\in} \Xi_{m, \beta_0}^{\text{lin}} \\ & \left( (\tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}'_i), \sum_{j \in [n^{\text{ip}}]} \mathbf{e}_{\iota_{\text{ip-in}}(j)} \otimes \hat{\mathbf{v}}_{j, i} \right) \stackrel{?}{\in} \Xi_{m, \beta_0}^{\text{lin}} \quad (i) \in [\ell] \end{aligned}$$

for the same parameters as in Fig. 8.

In sum, we obtain the reduction  $\Xi_{n^{\text{out}}, \beta}^{\text{eip}} \xrightarrow{II^{\text{eip}}} (\Xi_{n^{\text{out}}+3n^{\text{ip}}, \beta}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)}$  for the same parameters as the  $II^{\text{ip}}$  from Section 5.7. Next, we outline the full protocol  $II^{\text{eip}+}$  analogous to  $II^{\text{ip}+}$ .

**The full protocol** The full protocol  $II^{\text{eip}+}$  reduces the  $n^{\text{ip}} \cdot (\ell+1)$  statements after  $II^{\text{eip}}$  back to a single one with the minimal number of rows, by applying  $II^{\text{batch}}$  and  $II^{\text{fold}}$ . That is,  $II^{\text{eip}+}$  works as follows:

- (1)  $\Xi_{n^{\text{out}}, \beta}^{\text{ip}} \xrightarrow{II^{\text{eip}}} (\Xi_{n^{\text{out}}+3 \cdot n^{\text{ip}}, \beta}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)}$  (Reduce to claims in  $(\Xi^{\text{lin}})^{\ell+1}$ .)
- (2)  $(\Xi_{n^{\text{out}}+3 \cdot n^{\text{ip}}, \beta}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)} \xrightarrow{II^{\text{batch}}} (\Xi_{\bar{n}+1, \beta}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)}$  (Reduce number of rows again.)
- (3)  $(\Xi_{\bar{n}+1, \beta}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)} \xrightarrow{II^{\text{fold}}} \Xi_{\bar{n}+1, n^{\text{ip}} \cdot (\ell+1) \gamma \beta}^{\text{lin}}$  (Reduce to claim  $\beta' = (n^{\text{ip}} \cdot (\ell+1)) \gamma \beta$ .)

We do not provide explicit proof of Theorem 13 below as it completely analogous to Theorem 6. We denote by  $\Xi^{\text{eipVsis}}$  the usual relaxed relation.

**Theorem 13 (Security of  $II^{\text{eip}+}$ ).** *The protocol  $II^{\text{eip}+}$  is a reduction of knowledge  $\Xi_{m, n^{\text{out}}, \beta_0}^{\text{eip}}$  to  $(\Xi_{m, n^{\text{out}}+3 \cdot n^{\text{ip}}, \beta'_3 = n^{\text{ip}} \cdot (\ell+1) \beta_0}^{\text{lin}})^{n^{\text{ip}} \cdot (\ell+1)}$ , where  $\ell \geq \log_{b_{\text{ip}}} (2\beta^2 + 1)$  for  $b_{\text{ip}} \leq 2\beta / (\sqrt{m} \varphi^{3/2})$  for  $b_{\text{ip}}, \ell \in \mathbb{N}$ . It is perfectly correct. It is relaxed knowledge sound from  $\Xi_{m, n^{\text{out}}, \beta'_0}^{\text{eipVsis}}$  to  $\Xi_{m, n^{\text{out}}+3 \cdot n^{\text{ip}}, \beta'_3}^{\text{linVsis}}$  if  $\beta'_0 \leq \beta^{\text{sis}}$  where  $\beta'_0 = 4\theta\beta'_3$ . The knowledge error is  $n^{\text{ip}} \cdot (\ell+1) \cdot (|\mathcal{C}_{\mathcal{R}}|^{-1} + 4|\mathcal{C}_{\mathcal{R}_q}|^{-1}) + 2m|\mathcal{C}_{\mathcal{R}_q}|^{-1}$ .*

## 7.5 $II^{\text{aut}}$ : Automorphism Check

For our CRT-based protocol (Section 7), we need to efficiently check automorphism relations between different blocks of the witness. Formally, we define the automorphism check relation below.

$$\Xi_{m, n^{\text{out}}, \beta, \beta^{\text{sis}}}^{\text{autVsis}} := \left\{ \begin{aligned} & ((\iota_{\text{ip}}, \iota_{\text{ip-in}}, \iota_{\text{aut}}, \iota_{\text{aut-in}}, \iota_{\text{aut-out}}, \mathbf{H}, \mathbf{F}, \mathbf{y}, \mathbf{t}), \mathbf{w} = (\mathbf{w}_k)_{k \in [n^{\text{blk}}]}) : \\ & \mathbf{H} \in \mathcal{R}_q^{n^{\text{out}} \times n}; \mathbf{F} \in \mathcal{R}_q^{n \times (d^{\otimes \mu} \cdot n^{\text{blk}})} \subseteq \mathcal{R}_q^{n \times m}; \mathbf{y} \in \mathcal{R}_q^n; \mathbf{t} \in \mathcal{R}_q^{n^{\text{ip}}} \\ & \left\{ \begin{aligned} & \iota_{\text{ip}} \in \{1, -1\}^{[n^{\text{ip}}]}; \iota_{\text{ip-in}} \in [n^{\text{blk}}]^{[n^{\text{ip}}]} \\ & \iota_{\text{aut}} \in [\varphi]^{[n^{\text{aut}}]}; \iota_{\text{aut-in}}, \iota_{\text{aut-out}} \in [n^{\text{blk}}]^{[n^{\text{aut}}]} \\ & \|\mathbf{w}\| \leq \beta; \mathbf{H} \cdot \mathbf{F}\mathbf{w} = \mathbf{H}\mathbf{y} \pmod{q} \\ & \langle \mathbf{w}_{\iota_{\text{ip-in}}(k)}, \alpha_{\iota_{\text{ip}}(k)}(\mathbf{w}_{\iota_{\text{ip-in}}(k)}) \rangle = t_k \pmod{q} \quad \forall k \in [n^{\text{ip}}] \\ & \alpha_{\iota_{\text{aut}}(k)}(\mathbf{w}_{\iota_{\text{aut-in}}(k)}) = \mathbf{w}_{\iota_{\text{aut-out}}(k)} \pmod{q} \quad \forall k \in [n^{\text{aut}}] \end{aligned} \right\} \quad \text{or} \quad \left\{ \begin{aligned} & \|\mathbf{w}\| \leq \beta^{\text{sis}} \\ & \overline{\mathbf{H}}\mathbf{F}\mathbf{w} = \mathbf{0}_{\bar{n}} \pmod{q} \end{aligned} \right\} \end{aligned} \right\}.$$

Note that  $\Xi^{\text{aut}}$  implicitly has an additional parameter – the number of automorphism relations  $n^{\text{aut}}$ . An instance of  $\Xi^{\text{aut}}$  is almost identical to an instance of  $\Xi^{\text{eip}}$ , except that it additionally asserts that the automorphism relations  $\alpha_{\iota_{\text{aut}}(k)}(\mathbf{w}_{\iota_{\text{aut-in}}(k)}) = \mathbf{w}_{\iota_{\text{aut-out}}(k)} \bmod q$  hold for all  $k \in [n^{\text{aut}}]$  for the index maps  $\iota_{\text{aut}}, \iota_{\text{aut-in}}, \iota_{\text{aut-out}}$  given as part of the statement. Thus, to see that  $\Xi^{\text{aut}}$  reduces to  $\Xi^{\text{eip}}$ , we only need to check that the automorphism relations can be reduced to linear relations. We outline the logic of this reduction below, and give a formal description in Fig. 10.

Instead of checking that all automorphism relations hold, which would not be succinct, the verifier sends a random  $\xi \leftarrow \mathcal{C}_{\mathcal{R}_q} \subseteq \mathcal{R}_q$ . The prover then sends

$$\begin{aligned} z_k &:= (1, \xi, \dots, \xi^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-in}}(k)} \bmod q \\ z'_k &:= (1, \alpha_{\iota_{\text{aut}}(k)}(\xi), \dots, \alpha_{\iota_{\text{aut}}(k)}(\xi)^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-out}}(k)} \bmod q \end{aligned}$$

for  $k \in [n^{\text{aut}}]$ . In turn, the verifier checks that  $\alpha_{\iota_{\text{aut}}(k)}(z_k) = z'_k$  for  $k \in [n^{\text{aut}}]$ . Completeness can be seen by observing

$$\begin{aligned} \alpha_{\iota_{\text{aut}}(k)}(z_k) &= \alpha_{\iota_{\text{aut}}(k)}((1, \xi, \dots, \xi^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-in}}(k)}) \\ &= (1, \alpha_{\iota_{\text{aut}}(k)}(\xi), \dots, \alpha_{\iota_{\text{aut}}(k)}(\xi)^{m-1}) \cdot \alpha_{\iota_{\text{aut}}(k)}(\mathbf{w}_{\iota_{\text{aut-in}}(k)}) \\ &= (1, \alpha_{\iota_{\text{aut}}(k)}(\xi), \dots, \alpha_{\iota_{\text{aut}}(k)}(\xi)^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-out}}(k)} = z'_k. \end{aligned}$$

This reduces the automorphism check to verifying the linear relations

$$\begin{aligned} z_k &:= (1, \xi, \dots, \xi^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-in}}(k)} \bmod q \quad \forall k \in [n^{\text{aut}}], \\ z'_k &:= (1, \alpha_{\iota_{\text{aut}}(k)}(\xi), \dots, \alpha_{\iota_{\text{aut}}(k)}(\xi)^{m-1}) \cdot \mathbf{w}_{\iota_{\text{aut-out}}(k)} \bmod q \quad \forall k \in [n^{\text{aut}}] \end{aligned}$$

which is supported by  $\Xi^{\text{eip}}$ . Note that the reduction increases the number of rows significantly, and  $\Pi^{\text{batch}}$  could be run to batch them before further reductions.

The security of  $\Pi^{\text{aut}}$  is summarised below. The proof is omitted due to its similarity with the previous proofs in this section.

**Theorem 14 (Security of  $\Pi^{\text{aut}}$ ).** *The protocol  $\Pi^{\text{aut}}$  is a perfectly correct reduction of knowledge from  $\Xi_{m, n^{\text{out}}, \beta}^{\text{aut}}$  to  $\Xi_{m, n^{\text{out}}+2n^{\text{aut}}, \beta}^{\text{eip}}$ . It is relaxed knowledge sound from  $\Xi_{m, n^{\text{out}}, \beta'_0}^{\text{aut} \vee \text{sis}}$  to  $\Xi_{m, n^{\text{out}}+2n^{\text{aut}}, \beta'_1}^{\text{aut} \vee \text{sis}}$  if  $2\beta_1 \leq \beta^{\text{sis}}$ , where  $\beta'_0 = \beta'_1$ . The knowledge error is  $n^{\text{aut}}/|\mathcal{C}_{\mathcal{R}_q}|$ .*

## 7.6 Reducing CRT-based Binariness Check to Automorphism Check

Equipped with Theorems 12 and 14, in Fig. 11, we construct a reduction of knowledge  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$  using a CRT-based binariness check. For the ease of notation, the reduction of knowledge  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$  in Fig. 11 makes use of the following subroutine  $\text{CRTmake}(\mathbf{H}, \mathbf{F}, \mathbf{y}, \mathbf{w})$  which maps a  $\Xi^{\text{lin}} \cap \Xi^{\text{bin}}$  instance to a  $\Xi^{\text{aut}}$  instance:

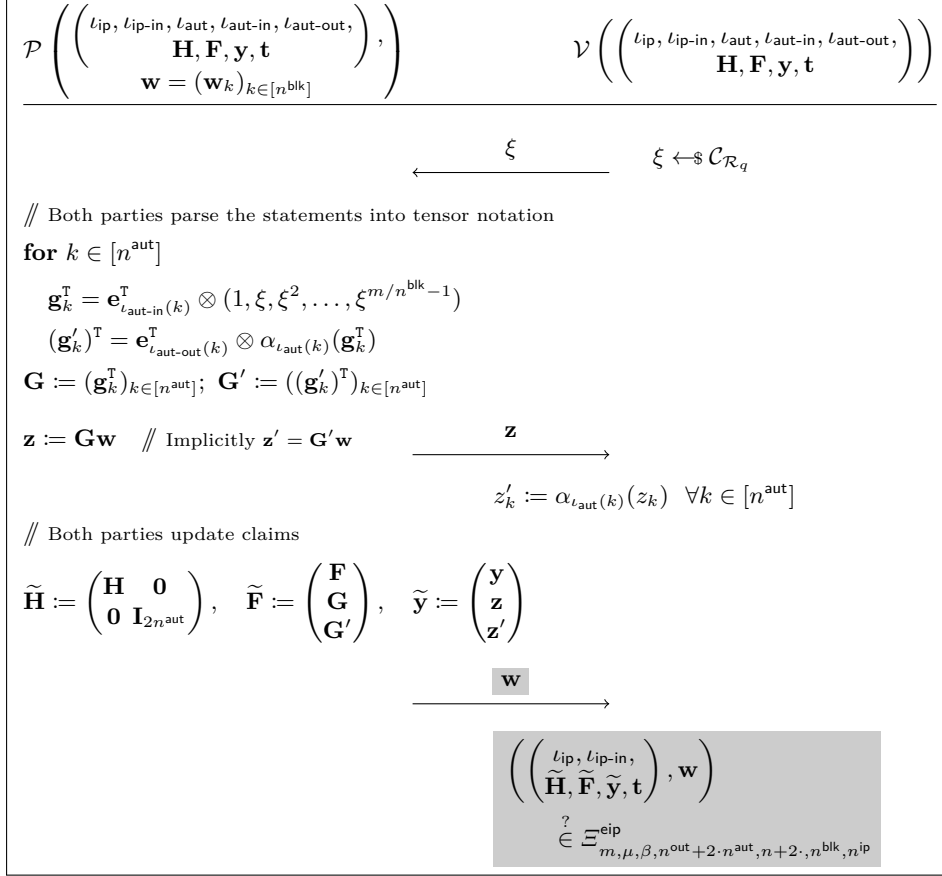
- (i) Compute  $\widehat{\mathbf{w}}_0 := \text{CRT}_p^{-1}(\psi(\mathbf{w}))$  and let  $\mathbf{r}_0 := \mathbf{0}$ .
- (ii) Compute  $\mathbf{y}' := \langle \widehat{\mathbf{w}}_0, \text{CRT}^{-1}(\mathbf{1}^{\delta m}) \rangle_{\mathcal{R}}$ ,  $\tilde{t}_0 := \langle \widehat{\mathbf{w}}_0, \widehat{\mathbf{w}}_0 \rangle_{\mathcal{R}}$ , and  $\tilde{t}_1 := \langle \mathbf{w}, \overline{\mathbf{w}} \rangle_{\mathcal{R}}$ .
- (iii) For each  $i \in [t]$ : (i) Compute the  $i$ -th step of the CRT decomposition as described in Theorem 12, i.e.  $\mathbf{w}_{i,j} := \alpha_{i,j}(\widehat{\mathbf{w}}_i)$  for all  $j \in [h_i]$ . (ii) Compute the  $p$ -ary representative  $\widehat{\mathbf{w}}_{i+1} := \sum_{j \in [h_i]} s_{i,j} \mathbf{w}_{i,j} \bmod p$ .
- (iii) Find the quotient  $\mathbf{r}_{i+1} \in \mathcal{R}^m$  satisfying  $\widehat{\mathbf{w}}_{i+1} - \sum_{j \in [h_i]} s_{i,j} \cdot \mathbf{w}_{i,j} + p \cdot \mathbf{r}_{i+1} = \mathbf{0}^m \pmod{q}$ .
- (iv) Concatenate all intermediate witnesses as

$$\widetilde{\mathbf{w}}^{\text{T}} := ((\mathbf{w}_{i,j}^{\text{T}})_{(i,j) \in [t, h_i]}, (\widehat{\mathbf{w}}_i^{\text{T}}, \mathbf{r}_i^{\text{T}})_{i \in [t+1]}) \in \mathcal{R}^{\widetilde{m}}$$

where  $\widetilde{m} := m \cdot n^{\text{blk}}$  and  $n^{\text{blk}} := w \log \mathfrak{f} + 2t + 2$ .<sup>19</sup>

- (v) Parse the witness as a block vector  $\widetilde{\mathbf{w}} = (\widetilde{\mathbf{w}}_k)_{k \in [n^{\text{blk}}]}$  where  $\widetilde{\mathbf{w}}_k \in \mathcal{R}^m$ .
- (vi) Concatenate all linear map images as  $\widetilde{\mathbf{y}}^{\text{T}} := (\mathbf{y}^{\text{T}}, \mathbf{0}^{tm}, c) \in \mathcal{R}_q^{\widetilde{n}}$  where  $\widetilde{n} := n + tm + 1$ .
- (vii) Concatenate all inner product images as  $\widetilde{\mathbf{t}}^{\text{T}} := (\tilde{t}_0, \tilde{t}_1) \in \mathcal{R}^2$ .

<sup>19</sup>Recall that  $\mathfrak{f}$  is the conductor of  $\mathcal{R}$  and is  $w$ -smooth. Technically,  $\widetilde{\mathbf{w}}$  has dimension  $m \cdot (\sum_{i \in [t]} h_i + 2t + 2)$ . Since  $h_i \leq w$  for all  $i$  and  $t \leq \log \mathfrak{f}$ , we pad  $\widetilde{\mathbf{w}}$  to dimension  $\widetilde{m} = m \cdot (w \log \mathfrak{f} + 2t + 2)$  for simplicity.



**Fig. 10.** Protocol  $\Pi^{\text{aut}}$ , a reduction of knowledge from  $\Xi_{m, \mu, \beta, n^{\text{out}}, n, n^{\text{blk}}, n^{\text{ip}}, n^{\text{aut}}}$  to  $\Xi_{m, \mu, \beta, n^{\text{out}}+2 \cdot n^{\text{aut}}, n+2 \cdot n^{\text{blk}}, n^{\text{ip}}}$ .  $\Pi^{\text{aut}}$  sends the marked parts only as a *proof* (but not *reduction*) of knowledge.

(viii) Define a matrix  $\tilde{\mathbf{F}} \in \mathcal{R}_q^{\tilde{n} \times \tilde{m}}$  such that  $\tilde{\mathbf{F}}\tilde{\mathbf{w}} = \tilde{\mathbf{y}} \pmod q$  represented the following system of equations:

$$\begin{aligned} \mathbf{F}\hat{\mathbf{w}}_t &= \mathbf{y} \pmod q, \\ \langle \text{CRT}^{-1}(\mathbf{1}^{\delta m}), \hat{\mathbf{w}}_0 \rangle_{\mathcal{R}} &= y' \pmod q, \\ \hat{\mathbf{w}}_{i+1} - \sum_{j \in [h_i]} s_{i,j} \cdot \mathbf{w}_{i,j} + p \cdot \mathbf{r}_{i+1} &= \mathbf{0}^m \pmod q \quad \forall i \in [t]. \end{aligned}$$

(ix) Write  $\{\sigma_k\}_{k \in \mathbb{Z}_f^\times} = \text{Gal}(\mathcal{K}/\mathbb{Q})$ .

(x) Let  $n^{\text{ip}} := 2$  and define the index maps  $\iota_{ip} : [n^{\text{ip}}] \rightarrow \{1, -1\}$  and  $\iota_{ip-in} : [n^{\text{ip}}] \rightarrow [n^{\text{blk}}]$  for inner products such that

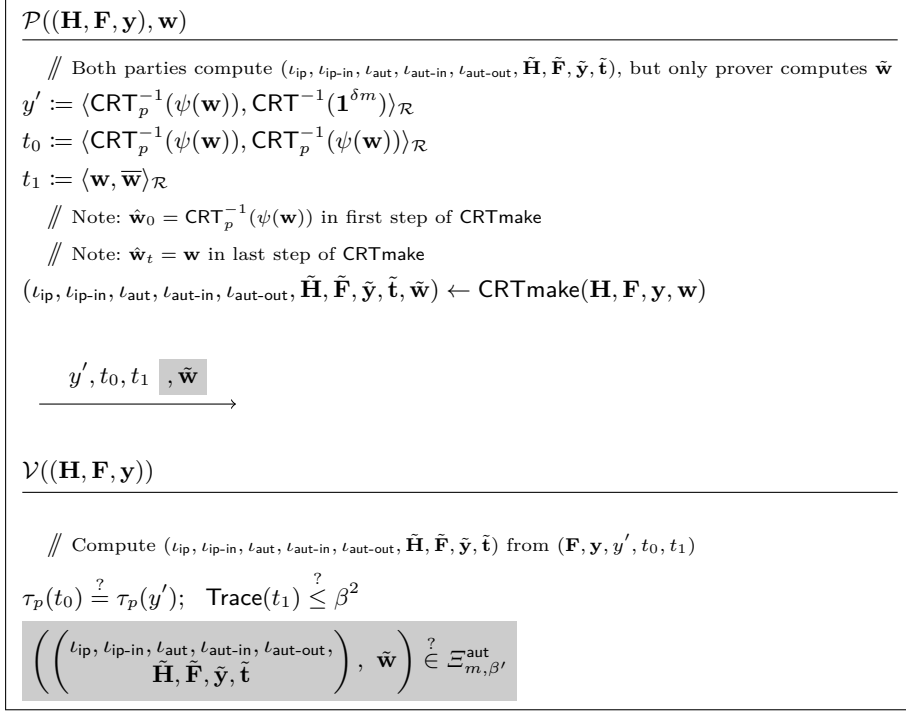
$$\forall k \in [n^{\text{ip}}], t_k = \langle \tilde{\mathbf{w}}_{\iota_{ip-in}(k)}, \sigma_{\iota_{ip}(k)}(\tilde{\mathbf{w}}_{\iota_{ip-in}(k)}) \rangle \iff \tilde{t}_0 = \langle \hat{\mathbf{w}}_0, \hat{\mathbf{w}}_0 \rangle_{\mathcal{R}} \quad \wedge \quad \tilde{t}_1 = \langle \hat{\mathbf{w}}_t, \overline{\hat{\mathbf{w}}_t} \rangle_{\mathcal{R}}$$

(xi) Let  $n^{\text{aut}} := w \log f$  and define the index maps  $\iota_{aut} : [n^{\text{aut}}] \rightarrow \mathbb{Z}_f^\times$  and  $\iota_{aut-in} : [n^{\text{aut}}] \rightarrow [n^{\text{blk}}]$ , and  $\iota_{aut-out} : [n^{\text{aut}}] \rightarrow [n^{\text{blk}}]$  for automorphisms so that

$$\forall k \in [f], \tilde{\mathbf{w}}_{\iota_{aut-out}(k)} = \sigma_{\iota_{aut}(k)}(\tilde{\mathbf{w}}_{\iota_{aut-in}(k)}) \iff \forall i \in [t], j \in [h_i], \mathbf{w}_{i,j} := \alpha_{i,j}(\hat{\mathbf{w}}_i).$$

(xii) Set  $\tilde{\mathbf{H}} := \begin{pmatrix} \mathbf{H} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{tm+1} \end{pmatrix}$ .

(xiii) Output  $(l_{ip}, l_{ip-in}, l_{aut}, l_{aut-in}, l_{aut-out}, \tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}, \tilde{\mathbf{t}}, \tilde{\mathbf{w}})$ .



**Fig. 11.** Protocol  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$ , a reduction from  $\Xi^{\text{lin}} \cap \Xi^{\text{bin}}$  to  $\Xi^{\text{aut}}$ , where  $\beta, \beta'$ , and  $p$  are set as in Theorem 15. See the text for the definition of the subroutine CRTmake. The marked parts are only sent / checked when  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$  is used as a proof of knowledge. As a reduction of knowledge, they are omitted.

**Theorem 15.** *The protocol  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$  is a perfectly correct from  $\Xi_{m, n^{\text{out}}, \beta}^{\text{lin}} \cap \Xi_{m, n^{\text{out}'}}^{\text{bin}}$  to  $\Xi_{m, n^{\text{out}}, \beta'}^{\text{aut}}$ , and knowledge sound reduction of knowledge from  $\Xi_{m, n^{\text{out}}, \beta}^{\text{lin}\vee\text{sis}} \cap \Xi_{m, n^{\text{out}'}}^{\text{bin}\vee\text{sis}}$  to  $\Xi_{m, n^{\text{out}}, \beta'}^{\text{aut}\vee\text{sis}}$ , where  $n^{\text{blk}} = w(1+t) + 2t + 2$ ,  $n^{\text{ip}} = 2$ ,  $n^{\text{aut}} = w \log \mathfrak{f}$ ,  $\beta' = \frac{p}{4} \cdot w \cdot \gamma_{\mathcal{R}} \cdot \varphi^{3/2} \cdot \sqrt{m} \cdot n^{\text{blk}}$ ,  $\beta = \varphi^{3/2} \sqrt{m}$ ,  $p/2 > c_{\mathfrak{f}}^2 \cdot \beta^2 \cdot m \cdot \varphi$ ,  $c_{\mathfrak{f}}$  is a constant from Corollary 2,  $n^{\text{out}'} = n^{\text{out}} + tm + 1$  and conductor  $\mathfrak{f}$  of the ring is  $w$ -smooth.*

*Proof. Completeness.* For perfect completeness, we consider the statement  $((\mathbf{H}, \mathbf{F}, \mathbf{y}), \mathbf{w}) \in \Xi_{m, \beta}^{\text{lin}} \cap \Xi_m^{\text{bin}}$ . Clearly, we have  $\psi(\mathbf{w}) \in \{0, 1\}^{m\delta}$ . Therefore,

$$\langle \psi(\mathbf{w}), \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} = \langle \psi(\mathbf{w}), \mathbf{w} \rangle_{\mathbb{Z}}$$

by Proposition 2 (and the discussion immediately after). Since  $\text{CRT}_p$  is a  $\tau_p$ -embedding over  $\mathcal{R}$  and  $\tilde{\mathbf{w}}_t = \mathbf{w}$  due to the construction it holds that

$$\begin{aligned} \tau_p(y') &= \langle \psi(\mathbf{w}), \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} \bmod p, \text{ and} \\ \tau_p(t_0) &= \langle \psi(\mathbf{w}), \psi(\mathbf{w}) \rangle_{\mathbb{Z}} \bmod p, \end{aligned}$$

As a consequence,  $\tau_p(t_0) = \tau_p(y')$  and the first check of the verifier passes.

Further, as  $\|\psi(\mathbf{w})\|_{\infty} \leq 1$ , and by Corollary 2  $\|\sigma(\mathbf{w})\|_2 \leq \varphi^{3/2} \sqrt{m} = \beta$ . We observe that  $\|\sigma(\mathbf{w})\|_2^2 = \text{Trace}(d)$  and hence the second verifier's check passes.

Next, we show that

$$\left( \left( \begin{array}{c} \iota_{ip}, \iota_{ip-in}, \iota_{aut}, \iota_{aut-in}, \iota_{aut-out}, \\ \tilde{\mathbf{H}}, \tilde{\mathbf{F}}, \tilde{\mathbf{y}}, \tilde{\mathbf{t}} \end{array} \right), \tilde{\mathbf{w}} \right) \in \Xi_{m, \mu, \beta', n^{\text{out}}, \tilde{n}, n^{\text{blk}}, n^{\text{ip}}, n^{\text{aut}}}^{\text{aut}}.$$

The linear part of the relation  $\Xi^{\text{aut}}$  above holds due to Item (viii) of the CRTmake subroutine. Similarly, the correctness of all steps of the CRT transformation implies that all automorphism relations hold. For the inner-product type of relation, they are correct due to the honest computation of  $t_0$  and  $t_1$ .

Finally, it remains to argue about the correctness of the bound  $\beta'$  for the new witness. Observe that the witness  $\tilde{\mathbf{w}}$  is a concatenation of three types of blocks. Particularly:

- the steps of the CRT transformation,  $\widehat{\mathbf{w}}_i$ ,
- the remainders,  $\mathbf{r}_i$ ,
- and  $\mathbf{w}_{i,j}$  used for automorphisms-based relations.

In the first case, due to reduction modulo  $p$ ,

$$\|\psi((\widehat{\mathbf{w}}_i)_{i \in [t+1]})\|_\infty \leq p/2.$$

After translating the norm, we consider the canonical 2-norm of each element of  $(\widehat{\mathbf{w}}_i)_{i \in [t+1]}$ .

$$\|\sigma(\widehat{w})\|_2 \leq p/2 \cdot \varphi^{3/2} \cdot \sqrt{m \cdot (1+t)} \quad \forall \widehat{w} \in (\widehat{\mathbf{w}}_i)_{i \in [t+1]}.$$

In the second case,

$$\|\sigma(w)\|_2 \leq p/2 \cdot \varphi^{3/2} \cdot \sqrt{m \cdot (1+t)} \quad \forall w \in (\mathbf{w}_{i,j})_{(i,j) \in [t,h_j]}.$$

as  $\mathbf{w}_{i,j} = \alpha_{i,j}(\widehat{\mathbf{w}}_i)$  and the automorphisms  $\alpha_{i,j}$  do not impact the canonical norm of a ring element.

In the third case,

$$\|\psi(r)\|_\infty \leq \left(\frac{p}{2} \cdot \frac{p}{2} \cdot w \cdot \gamma_{\mathcal{R}}\right) / p = \frac{p}{4} \cdot w \cdot \gamma_{\mathcal{R}} \quad \forall r \in (\mathbf{r}_i)_{i \in [t+1]}$$

and after translating norms,

$$\|\sigma(r)\|_2 \leq \frac{p}{4} \cdot w \cdot \gamma_{\mathcal{R}} \cdot \varphi^{3/2} \quad \forall r \in (\mathbf{r}_i)_{i \in [t+1]}.$$

To sum up, canonical 2-norm of individual elements of  $\widetilde{\mathbf{w}}$  is upper-bounded by

$$\|\sigma(\widetilde{w})\|_2 \leq \frac{p}{4} \cdot w \cdot \gamma_{\mathcal{R}} \cdot \varphi^{3/2} \quad \forall \widetilde{w} \in \widetilde{\mathbf{w}}$$

as  $w \geq 2$ . Last, we observe that  $|(\widehat{\mathbf{w}}_i)_{i \in [t+1]}| = m \cdot (t+1)$ ,  $|(\mathbf{w}_{i,j})_{(i,j) \in [t,h_j]}| \leq m \cdot (t+1) \cdot w$  and  $|(\mathbf{r}_i)_{i \in [t+1]}| = m \cdot (t+1)$ . As a result  $|\widetilde{\mathbf{w}}| \leq m \cdot (w(1+t) + 2t + 2) = n^{\text{blk}}$  so we derive

$$\|\sigma(\widetilde{\mathbf{w}})\|_2 \leq \frac{p}{4} \cdot w \cdot \gamma_{\mathcal{R}} \cdot \varphi^{3/2} \cdot \sqrt{m \cdot n^{\text{blk}}} = \beta',$$

which concludes the proof of the perfect correctness.

*Soundness.* For perfect knowledge soundness, suppose that

$$\left( \left( \iota_{\text{ip}}, \iota_{\text{ip-in}}, \iota_{\text{aut}}, \iota_{\text{aut-in}}, \iota_{\text{aut-out}}, \widetilde{\mathbf{H}}, \widetilde{\mathbf{F}}, \widetilde{\mathbf{y}}, \widetilde{\mathbf{t}} \right), \widetilde{\mathbf{w}} \right) \in \Xi_{m, \mu, \beta', n^{\text{out}}, \widetilde{n}, n^{\text{blk}}, n^{\text{ip}}, n^{\text{aut}}}^{\text{aut}}$$

meaning that  $\left( \iota_{\text{ip}}, \iota_{\text{ip-in}}, \iota_{\text{aut}}, \iota_{\text{aut-in}}, \iota_{\text{aut-out}}, \widetilde{\mathbf{H}}, \widetilde{\mathbf{F}}, \widetilde{\mathbf{y}}, \widetilde{\mathbf{t}} \right)$  takes the form as constructed in CRTmake. We have

$$\tau_p(y') = \langle \psi(\mathbf{w}), \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} \bmod p \text{ and } \tau_p(t_0) = \langle \psi(\mathbf{w}), \psi(\mathbf{w}) \rangle_{\mathbb{Z}} \bmod p.$$

Since  $\tau_p(y') = \tau_p(t_0)$ , we have

$$\langle \psi(\mathbf{w}), \psi(\mathbf{w}) - \mathbf{1}^{m\delta} \rangle_{\mathbb{Z}} = 0 \bmod p.$$

Furthermore,  $\text{Trace}(t_1) \leq \beta^2$  implies that  $\|\sigma(\mathbf{w})\|_2 \leq \beta$ , thus by Corollary 2  $\|\psi(\mathbf{w})\|_\infty \leq c_f \cdot \beta$ , where  $c_f$  is dependent only on  $f$ . Therefore, as  $c_f^2 \cdot \beta^2 \cdot m \cdot \varphi < p/2$ , the inner product above holds without modulo  $p$ , and thus  $\psi(\mathbf{w}) \in \{0, 1\}^{m\delta}$ .  $\square$

*Remark 8.* Theorem 15 trivially generalises to yield a reduction of knowledge from  $\Xi^{\text{ip}} \cap \Xi^{\text{bin}}$  to  $\Xi^{\text{aut}}$ , analogous to Theorem 11. We omit the details for succinctness.

*Remark 9.* Note that the reduction of knowledge  $\Pi_{\text{crt}, p}^{\text{lin-bin}}$  increases the number of linear relations from  $n$  to  $n + tm + 1$ , where we recall that  $m$  is the dimension for each block of the witness. Nevertheless, we observe that the matrices  $\widetilde{\mathbf{H}}$  and  $\widetilde{\mathbf{F}}$  in the statement are sparse and highly structured. Therefore, when applying the chain of reductions in Section 5, the verifier time can still be polylogarithmic in  $m$  as long as succinct representations of  $\widetilde{\mathbf{H}}$  and  $\widetilde{\mathbf{F}}$  are used when running  $\Pi^{\text{batch}}$  for batching.

*Remark 10.* The CRT-based embedding method discussed in this section requires the prime  $p$  to split completely, i.e.  $p = 1 \pmod{\mathfrak{f}}$ , which immediately implies that  $\mathfrak{f}$  must be even. When  $\mathfrak{f} = 2^k$ , we observe that choices of  $p$  are limited, and they tend to be significantly larger than  $\mathfrak{f}$  – the smallest values of  $p$  for  $\mathfrak{f} \in (256, 512, 1024)$  are  $p \in (267, 7681, 12289)$ . For general  $\mathfrak{f}$ , choices appear to be more flexible – the smallest values of  $p$  for  $\mathfrak{f} \in (598, 1102, 2926)$  are  $p \in (599, 1103, 2927)$ .

*Remark 11.* We discuss a possible optimisation which allows to pick smaller primes  $p$ . Recall that, for knowledge soundness, we required  $p/2 > c_{\mathfrak{f}}^2 \cdot \beta^2 \cdot m \cdot \varphi$ , which makes  $p$  linear in the length of the witness length. Towards picking smaller primes  $p$ , we suggest using multiple primes  $p_1, p_2, \dots$  such that all of which fully split over  $\mathcal{R}$ . To be concrete, consider the case with two primes, i.e.  $p_1$  and  $p_2$ . After repeating the protocol twice for  $p_1$  and  $p_2$  respectively (or even better, running a merged version of the protocol), the verifier should be convinced that:

$$\left. \begin{array}{l} \langle \mathbf{x}, \mathbf{y} \rangle = c \pmod{p_1} \\ \langle \mathbf{x}, \mathbf{y} \rangle = c \pmod{p_2} \\ (p_1, p_2) = 1 \end{array} \right\} \implies \langle \mathbf{x}, \mathbf{y} \rangle = c \pmod{p_1 \cdot p_2}.$$

Obviously, this generalises to having a set  $P$  of arbitrarily many primes. Consequently, for knowledge soundness, we would only require  $\prod_{p \in P} p/2 > c_{\mathfrak{f}}^2 \cdot \beta^2 \cdot m \cdot \varphi$ . If the conductor  $\mathfrak{f}$  is smooth, it is possible to find many highly favourable primes (cf. Remark 10).

## 8 Parameter Selection

We propose concrete instantiations of our protocols for various values of  $m$ . For comparison with prior works, e.g. [BS23, AFLN23], we aim for 128-bit security. This corresponds to the root Hermite factor  $\delta_{\text{rhf}} \approx 1.0044$  (cf. Section 5.9).

### 8.1 Split-and-Fold with Norm Checks

We start with instantiating the split-and-fold with an intermediate norm check described in Section 5.8. We focus on the following simple goal: commit to a short vector  $\mathbf{w} \in \mathbb{Z}_q^h$  of length  $h$ , such that  $\|\psi(\mathbf{w})\|_{\infty} \leq \beta_{\text{init}} = 2^5$  and prove knowledge of the commitment opening. To this end, we will pack  $h = m \cdot \varphi$  integers into a vector  $\mathbf{w} \in \mathcal{R}_q^m$  of  $m$  ring elements employing standard coefficient embedding. Then, we will use the vSIS commitment scheme on  $\mathbf{w}$ .

The relation of our interest is a proof of vSIS commitment opening [CLM23], i.e. the polynomial evaluation equation  $\mathbf{w}(v) = y \pmod{q}$  for public ring elements  $v, y \in \mathcal{R}_q$ . When adapting this relation to the language of  $\Xi^{\text{ip}}$ , we would initially set  $(n, n^{\text{out}}) = (1, 1)$ . Throughout the batching protocols, we set  $\bar{n} = 1$ . In our experiments, we hardwire  $\mathfrak{f} = 5544$ . Hence,  $\varphi = 1440$  and  $|\mathcal{C}_{\mathcal{R}}| = 504$ . We fix  $\ell = 2$  for the whole execution of the protocol. Then, given the norm bound  $\beta$  and  $\ell$ , we can deduce the decomposition base  $b$ . Since in each iteration of the split-and-fold protocol, the norm  $\beta$  may change, then so can the base  $b$ . We note that (at least concretely) the current proof sizes are not optimal, reaching high orders of Megabytes. We highlight that we could achieve better sizes for larger moduli. However, conceptually this would be contradictory to the original intention of split-and-fold with norm checks; avoiding unnecessary stretch and large proof system modulus. We defer more fine-grained optimisation to a follow-up work.

**Proof Composition.** To further shrink communication, one could use standard proof composition, where instead of the verifier checking the verification conditions, the prover provides a proof of knowledge of the input for which the verification holds. To this end, we can directly apply the LaBRADOR proof system [BS23]. Note that this approach is different from running [BS23] for the original relation because now the statement/witness size for LaBRADOR is only of size  $\text{poly}(\lambda, \log m)$ , and thus we do maintain succinct verification. Hence, we estimate the final communication size to be  $\approx 100\text{KB}$  based on performance described in [BS23].



witness length in $\mathbb{Z}$ -elements	$\approx 2^{18}$	$\approx 2^{20}$	$\approx 2^{24}$
$\log q$	110	110	120
$f$	5544	5544	5544
$m$	$2^{20}$	$2^{22}$	$2^{26}$
$\ \psi(\cdot)\ _\infty$ of the witness	$2^5$	$2^5$	$2^5$
witness size	1080 MB	4320 MB	69120 MB
$d$	2	2	2
$\mu$	2	3	4
# of repetitions	17	18	19
(unoptimized) proof size	258.4 MB	263.6 MB	458.6 MB

**Table 2.** Concrete parameters, together with proof sizes, for security level  $\lambda = 128$ .

**Fiat-Shamir Transformation.** As noted in [AFK22], transforming interactive proofs, which admit parallel repetition, to the non-interactive setting via Fiat-Shamir transformation incurs significant loss in the order of  $Q^\mu$ , where  $Q$  is the number of random oracle queries made by an adversary. In the asymptotic sense, here we could simply set  $d = O(\lambda)$ , and thus  $\mu = O(\frac{\log m}{\log \lambda})$  which is constant for  $m = \text{poly}(\lambda)$ . Unfortunately, as seen in the experiments, the most optimal proof sizes are achieved for small values of  $d$ . In this section, we heuristically omit the security loss  $Q^\mu$  and leave any potential improvements as future work.

**Concrete parameters.** In Table 2 we suggest concrete parameters with the estimated proof size. The results are obtained via a dedicated script<sup>20</sup> simulating protocol execution and measuring the communication cost.

## Acknowledgments

The work of R.L. and M.O. was supported by the Research Council of Finland project No. 358951. This work was supported by the Helsinki Institute for Information Technology (HIIT) and conducted while M.K. was affiliated with Aalto University. N.K.N. was supported by the Protocol Labs RFP-013: Cryptonet network grant.

## References

- ACL<sup>+</sup>22. Martin R. Albrecht, Valerio Cini, Russell W. F. Lai, Giulio Malavolta, and Sri Aravinda Krishnan Thyagarajan. Lattice-based SNARKs: Publicly verifiable, preprocessing, and recursively composable - (extended abstract). In Yevgeniy Dodis and Thomas Shrimpton, editors, *CRYPTO 2022, Part II*, volume 13508 of *LNCS*, pages 102–132. Springer, Heidelberg, August 2022. 1, 45
- AFK22. Thomas Attema, Serge Fehr, and Michael Klooß. Fiat-shamir transformation of multi-round interactive proofs. In Eike Kiltz and Vinod Vaikuntanathan, editors, *TCC 2022, Part I*, volume 13747 of *LNCS*, pages 113–142. Springer, Heidelberg, November 2022. 41
- AFLN23. Martin R. Albrecht, Giacomo Fenzi, Oleksandra Lapiha, and Ngoc Khanh Nguyen. Slap: Succinct lattice-based polynomial commitments from standard assumptions. Cryptology ePrint Archive, Paper 2023/1469, 2023. <https://eprint.iacr.org/2023/1469>. 1, 2, 4, 40
- AL21. Martin R. Albrecht and Russell W. F. Lai. Subtractive sets over cyclotomic rings - limits of Schnorr-like arguments over lattices. In Tal Malkin and Chris Peikert, editors, *CRYPTO 2021, Part II*, volume 12826 of *LNCS*, pages 519–548, Virtual Event, August 2021. Springer, Heidelberg. 2, 5, 10, 11
- BC24. Dan Boneh and Binyi Chen. Latticefold: A lattice-based folding scheme and its applications to succinct proof systems. Cryptology ePrint Archive, Paper 2024/257, 2024. <https://eprint.iacr.org/2024/257>. 3
- BCS16. Eli Ben-Sasson, Alessandro Chiesa, and Nicholas Spooner. Interactive oracle proofs. In Martin Hirt and Adam D. Smith, editors, *TCC 2016-B, Part II*, volume 9986 of *LNCS*, pages 31–60. Springer, Heidelberg, October / November 2016. 1

<sup>20</sup>The script and the output are available at <https://github.com/russell-lai/rok-paper-sissors-estimator/blob/camera-ready/rok-estimator.ipynb>

- BCS23. Jonathan Bootle, Alessandro Chiesa, and Katerina Sotiraki. Lattice-based succinct arguments for NP with polylogarithmic-time verification. In Helena Handschuh and Anna Lysyanskaya, editors, *CRYPTO 2023, Part II*, volume 14082 of *LNCS*, pages 227–251. Springer, Heidelberg, August 2023. [2](#), [4](#)
- BDGL15. Anja Becker, Léo Ducas, Nicolas Gama, and Thijs Laarhoven. New directions in nearest neighbor searching with applications to lattice sieving. Cryptology ePrint Archive, Report 2015/1128, 2015. <https://eprint.iacr.org/2015/1128>. [25](#)
- BFOV04. E. Bayer-Fluckiger, F. Oggier, and E. Viterbo. New algebraic constructions of rotated  $z$ -sup  $n$ -lattice constellations for the rayleigh fading channel. *IEEE Transactions on Information Theory*, 50(4):702–714, 2004. [8](#), [27](#), [28](#)
- BL17. Carsten Baum and Vadim Lyubashevsky. Simple amortized proofs of shortness for linear relations over polynomial rings. Cryptology ePrint Archive, Report 2017/759, 2017. <https://eprint.iacr.org/2017/759>. [2](#)
- BLNS20. Jonathan Bootle, Vadim Lyubashevsky, Ngoc Khanh Nguyen, and Gregor Seiler. A non-PCP approach to succinct quantum-safe zero-knowledge. In Daniele Micciancio and Thomas Ristenpart, editors, *CRYPTO 2020, Part II*, volume 12171 of *LNCS*, pages 441–469. Springer, Heidelberg, August 2020. [2](#)
- BS23. Ward Beullens and Gregor Seiler. LaBRADOR: Compact proofs for R1CS from module-SIS. In Helena Handschuh and Anna Lysyanskaya, editors, *CRYPTO 2023, Part V*, volume 14085 of *LNCS*, pages 518–548. Springer, Heidelberg, August 2023. [1](#), [2](#), [8](#), [32](#), [40](#)
- CLM23. Valerio Cini, Russell W. F. Lai, and Giulio Malavolta. Lattice-based succinct arguments from vanishing polynomials - (extended abstract). In Helena Handschuh and Anna Lysyanskaya, editors, *CRYPTO 2023, Part II*, volume 14082 of *LNCS*, pages 72–105. Springer, Heidelberg, August 2023. [1](#), [3](#), [4](#), [5](#), [9](#), [14](#), [40](#), [45](#)
- CMNW24. Valerio Cini, Giulio Malavolta, Ngoc Khanh Nguyen, and Hoeteck Wee. Polynomial commitments from lattices: Post-quantum security, fast verification and transparent setup. Cryptology ePrint Archive, Paper 2024/281, 2024. <https://eprint.iacr.org/2024/281>. [1](#), [4](#)
- DAFS24. Thomas Debris-Alazard, Pouria Fallahpour, and Damien Stehlé. Quantum oblivious lwe sampling and insecurity of standard model lattice-based snarks. Cryptology ePrint Archive, Paper 2024/030, 2024. <https://eprint.iacr.org/2024/030>. [1](#)
- DKL<sup>+</sup>18. Léo Ducas, Eike Kiltz, Tancrede Lepoint, Vadim Lyubashevsky, Peter Schwabe, Gregor Seiler, and Damien Stehlé. Crystals-dilithium: A lattice-based digital signature scheme. *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2018(1):238–268, Feb. 2018. [2](#)
- DPSZ11. I. Damgård, V. Pastro, N.P. Smart, and S. Zakarias. Multiparty computation from somewhat homomorphic encryption. Cryptology ePrint Archive, Report 2011/535, 2011. <https://eprint.iacr.org/2011/535>. [31](#), [44](#)
- DPSZ12. Ivan Damgård, Valerio Pastro, Nigel P. Smart, and Sarah Zakarias. Multiparty computation from somewhat homomorphic encryption. In Reihaneh Safavi-Naini and Ran Canetti, editors, *CRYPTO 2012*, volume 7417 of *LNCS*, pages 643–662. Springer, Heidelberg, August 2012. [44](#)
- EKS<sup>+</sup>21. Muhammed F. Esgin, Veronika Kuchta, Amin Sakzad, Ron Steinfeld, Zhenfei Zhang, Shifeng Sun, and Shumo Chu. Practical post-quantum few-time verifiable random function with applications to algorand. In Nikita Borisov and Claudia Díaz, editors, *FC 2021, Part II*, volume 12675 of *LNCS*, pages 560–578. Springer, Heidelberg, March 2021. [2](#)
- FMN23. Giacomo Fenzi, Hossein Moghaddas, and Ngoc Khanh Nguyen. Lattice-based polynomial commitments: Towards asymptotic and concrete efficiency. Cryptology ePrint Archive, Paper 2023/846, 2023. <https://eprint.iacr.org/2023/846>. [1](#), [4](#), [10](#), [45](#), [46](#)
- GHL22. Craig Gentry, Shai Halevi, and Vadim Lyubashevsky. Practical non-interactive publicly verifiable secret sharing with thousands of parties. In Orr Dunkelman and Stefan Dziembowski, editors, *EUROCRYPT 2022, Part I*, volume 13275 of *LNCS*, pages 458–487. Springer, Heidelberg, May / June 2022. [2](#)
- HKR19. Max Hoffmann, Michael Klooß, and Andy Rupp. Efficient zero-knowledge arguments in the discrete log setting, revisited. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *ACM CCS 2019*, pages 2093–2110. ACM Press, Heidelberg, November 2019. [6](#)
- KP23. Abhiram Kothapalli and Bryan Parno. Algebraic reductions of knowledge. In Helena Handschuh and Anna Lysyanskaya, editors, *CRYPTO 2023, Part IV*, volume 14084 of *LNCS*, pages 669–701. Springer, Heidelberg, August 2023. [9](#), [10](#), [45](#), [46](#)
- Len76. H. W. Lenstra. Euclidean number fields of large degree. *Inventiones mathematicae*, 38(3):237–254, 1976. [2](#), [5](#), [10](#), [11](#)
- LM23. Russell W. F. Lai and Giulio Malavolta. Lattice-based timed cryptography. In Helena Handschuh and Anna Lysyanskaya, editors, *CRYPTO 2023, Part V*, volume 14085 of *LNCS*, pages 782–804. Springer, Heidelberg, August 2023. [7](#)

- LN17. Vadim Lyubashevsky and Gregory Neven. One-shot verifiable encryption from lattices. In Jean-Sébastien Coron and Jesper Buus Nielsen, editors, *EUROCRYPT 2017, Part I*, volume 10210 of *LNCS*, pages 293–323. Springer, Heidelberg, April / May 2017. [2](#)
- LNP22. Vadim Lyubashevsky, Ngoc Khanh Nguyen, and Maxime Plançon. Lattice-based zero-knowledge proofs and applications: Shorter, simpler, and more general. In Yevgeniy Dodis and Thomas Shrimpton, editors, *CRYPTO 2022, Part II*, volume 13508 of *LNCS*, pages 71–101. Springer, Heidelberg, August 2022. [1](#), [2](#), [3](#), [26](#), [29](#)
- LNS20. Vadim Lyubashevsky, Ngoc Khanh Nguyen, and Gregor Seiler. Practical lattice-based zero-knowledge proofs for integer relations. In Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna, editors, *ACM CCS 2020*, pages 1051–1070. ACM Press, November 2020. [8](#), [32](#)
- LNS21. Vadim Lyubashevsky, Ngoc Khanh Nguyen, and Gregor Seiler. Shorter lattice-based zero-knowledge proofs via one-time commitments. In Juan Garay, editor, *PKC 2021, Part I*, volume 12710 of *LNCS*, pages 215–241. Springer, Heidelberg, May 2021. [2](#)
- LPR13. Vadim Lyubashevsky, Chris Peikert, and Oded Regev. A toolkit for ring-LWE cryptography. In Thomas Johansson and Phong Q. Nguyen, editors, *EUROCRYPT 2013*, volume 7881 of *LNCS*, pages 35–54. Springer, Heidelberg, May 2013. [44](#)
- Lyu12. Vadim Lyubashevsky. Lattice signatures without trapdoors. In David Pointcheval and Thomas Johansson, editors, *EUROCRYPT 2012*, volume 7237 of *LNCS*, pages 738–755. Springer, Heidelberg, April 2012. [1](#), [2](#)
- MR09. Daniele Micciancio and Oded Regev. *Lattice-based Cryptography*, pages 147–191. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009. [4](#), [25](#)
- PSTY13. Charalampos Papamanthou, Elaine Shi, Roberto Tamassia, and Ke Yi. Streaming authenticated data structures. In Thomas Johansson and Phong Q. Nguyen, editors, *EUROCRYPT 2013*, volume 7881 of *LNCS*, pages 353–370. Springer, Heidelberg, May 2013. [3](#)
- Was97. Lawrence C. Washington. *Introduction to cyclotomic fields*, volume 83. Springer, 1997. [12](#)
- WW23. Hoeteck Wee and David J. Wu. Lattice-based functional commitments: Fast verification and cryptanalysis. In Jian Guo and Ron Steinfeld, editors, *ASIACRYPT 2023, Part V*, volume 14442 of *LNCS*, pages 201–235. Springer, Heidelberg, December 2023. [1](#)

## A Extended Preliminaries

### A.1 Algebraic Number Theory

For a modulus  $q \in \mathbb{N}$ , we write  $\mathcal{R}_q := \mathcal{R}/q\mathcal{R}$ . We denote by  $\mathcal{R}^\times$  and  $\mathcal{R}_q^\times$  the sets of units in  $\mathcal{R}$  and  $\mathcal{R}_q$  respectively. We endow  $\mathcal{R}$  with two geometries via the coefficient embedding  $\psi_{\mathbf{b}} : \mathcal{R} \rightarrow \mathbb{Z}^\delta$  (for a given basis  $\mathbf{b}$ ) and the canonical embedding  $\sigma : \mathcal{K} \rightarrow \mathbb{C}^\varphi$  (of  $\mathcal{K}$ ). Specifically, for a given  $\mathbb{Z}$ -basis  $\mathbf{b} = (b_i)_{i \in [\delta]}$  of  $\mathcal{R}$  and an element  $x = \sum_{i \in [\delta]} x_i b_i \in \mathcal{R}$ , we write

$$\psi_{\mathbf{b}}(x) := (x_i)_{i \in [\delta]} \quad \text{and} \quad \sigma(x) := (\sigma_j(x))_{j \in [\varphi]}$$

where  $\sigma_j \in \text{Gal}(\mathcal{K}/\mathbb{Q})$ . Note that we define  $\sigma(x)$  by treating  $x \in \mathcal{K}$  in order to avoid discussing the canonical embedding of subfields of  $\mathcal{K}$ . If  $\mathcal{R} = \mathcal{O}_{\mathcal{K}}$  and  $\mathbf{b} = (1, \zeta, \dots, \zeta^{\varphi-1})$  is the standard power basis, we may omit  $\mathbf{b}$  from the subscript of  $\psi_{\mathbf{b}}$ . We extend the notation of  $\psi_{\mathbf{b}}$  and  $\sigma$  naturally to vectors, i.e. if  $\mathbf{x} = (x_i)_{i \in [m]} \in \mathcal{R}^m$ , then

$$\psi_{\mathbf{b}}(\mathbf{x}) := (\psi_{\mathbf{b}}(x_i))_{i \in [\delta]} \quad \text{and} \quad \sigma(\mathbf{x}) := (\sigma_j(x_i))_{j \in [\varphi]}$$

are defined as concatenations.

For any  $p \in \mathbb{N}$ , we consider the balanced representation of  $\mathbb{Z}_p$ , i.e. elements are represented by  $[-p/2, p/2) \cap \mathbb{Z}$ . When considering the quotient ring  $\mathcal{R}_p := \mathcal{R}/p\mathcal{R}$  where  $\mathcal{R}$  has  $\mathbb{Z}$ -basis  $\mathbf{b}$ , we assume that an element  $x \in \mathcal{R}_p$  is represented by  $\psi_{\mathbf{b}}(x) \in ([-p/2, p/2) \cap \mathbb{Z})^\delta$ . As such, for any  $x \in \mathcal{R}$ , we abuse the notation  $x \in \mathcal{R}_p$  to mean that  $\psi(x) \in ([-p/2, p/2) \cap \mathbb{Z})^\delta$ . The above extends naturally to vectors over  $\mathcal{R}$ .

To distinguish between  $\mathbb{Z}$ -inner products and  $\mathcal{R}$ -inner products, we write  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{Z}} = \sum_{i \in [m]} x_i y_i$  or  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{R}} = \sum_{i \in [m]} x_i y_i$  depending on whether  $\mathbf{x}, \mathbf{y} \in \mathbb{Z}^m$  or  $\mathbf{x}, \mathbf{y} \in \mathcal{R}^m$ . Note that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{R}}$  is defined without complex conjugation.

For any Galois extension  $\mathcal{M}/\mathcal{L}$ , the field trace can be computed as  $\text{Trace}_{\mathcal{M}/\mathcal{L}} : \mathcal{K} \rightarrow \mathcal{L}$ ,  $\text{Trace}_{\mathcal{M}/\mathcal{L}}(x) := \sum_{\sigma_j \in \text{Gal}(\mathcal{K}/\mathcal{L})} \sigma_j(x)$ . When  $\mathcal{L} = \mathbb{Q}$ , we drop the subscript and write  $\text{Trace} = \text{Trace}_{\mathcal{M}/\mathbb{Q}}$ .

The coefficient  $\ell_p$ -norm and canonical  $\ell_p$ -norm of a vector  $\mathbf{x} \in \mathcal{R}^m$  is denoted by  $\|\psi(\mathbf{x})\|_p$  and  $\|\sigma(\mathbf{x})\|_p$  respectively. We will mostly use  $\|\psi(\cdot)\|_\infty$  and  $\|\sigma(\cdot)\|_2$ . The ring expansion factor of  $\mathcal{R}$  w.r.t. the

coefficient  $\ell_\infty$ -norm is defined as  $\gamma_{\mathcal{R}} := \max_{a,b \in \mathcal{R}} \|\psi(a \cdot b)\|_\infty / (\|\psi(a)\|_\infty \cdot \|\psi(b)\|_\infty)$ . Assuming balanced representation, for any  $x \in \mathcal{R}_p$ , we have  $\|\psi(x)\|_\infty \leq p/2$ . Note that  $\|\sigma(\mathbf{x})\|_2^2 = \text{Trace}(\mathbf{x}^T \bar{\mathbf{x}})$ , where  $\bar{\cdot}$  denotes the complex conjugate.

**Lemma 9.** *Let  $x \in \mathcal{K}$ . It holds that  $\|\sigma(x)\|_2 \leq \varphi^{3/2} \|\psi(x)\|_\infty$*

*Proof.* Let  $x = \sum_{i \in [\varphi]} x_i \zeta^i$ , with  $x_i \in \mathbb{Q}$ . Then we have  $\|\sigma(x)\|_2 = \left\| \sigma \left( \sum_{i \in [\varphi]} x_i \zeta^i \right) \right\|_2 \leq \sum_{i \in [\varphi]} x_i \|\sigma(\zeta^i)\|_2 \leq \varphi \cdot \max_i (|x_i|) \cdot \sqrt{\varphi} = \varphi^{3/2} \|\psi(x)\|_\infty$ .

**Lemma 10 ([DPSZ11, Theorem 7], [DPSZ12]).** *For any conductor  $\mathfrak{f}$ , there exists a constant  $c_{\mathfrak{f}}$ , such that for any  $x \in \mathcal{K}$ , it holds that  $\|\psi(x)\|_\infty \leq c_{\mathfrak{f}} \cdot \|\sigma(x)\|_\infty$ .*

**Corollary 2.** *For any conductor  $\mathfrak{f}$ , there exists a constant  $c_{\mathfrak{f}}$ , such that for any  $x \in \mathcal{K}$ , it holds that  $\|\psi(x)\|_\infty \leq c_{\mathfrak{f}} \cdot \|\sigma(x)\|_2$ . More generally, for any  $\mathbf{x} \in \mathcal{K}^m$ , it holds that  $\|\psi(\mathbf{x})\|_\infty \leq c_{\mathfrak{f}} \cdot \|\sigma(\mathbf{x})\|_2$ .*

In [DPSZ11, Appendix C.2], it was discussed that  $c_{\mathfrak{f}} \leq 8.6$  for any square-free  $\mathfrak{f} \leq 400$ . However, note that in Lemma 10 the second inequality is trivial and usually not tight. If we only concern about bounding  $\|\psi(x)\|_\infty$  in terms of  $\|\sigma(x)\|_2$ , it seems that  $\|\psi(x)\|_\infty \leq \|\sigma(x)\|_2$  always holds. We verified this empirically for  $100 \leq \mathfrak{f} \leq 500$  with 100 random elements for each conductor.

**Lemma 11.** *If  $L/K$  is a Galois extension, then any  $K$ -linear map  $f : L \rightarrow L$  can be expressed as an  $L$ -linear combination of  $\text{Gal}(L/K)$ .*

*Proof.* Let  $L/K$  be of degree  $\varphi$ . Let  $\mathbf{e} = (e_i)_{i \in [\varphi]}$  be a  $\mathcal{O}_K$ -basis of  $\mathcal{O}_L$ , and  $\mathbf{e}^\vee = (e_i^\vee)_{i \in [\varphi]}$  be a  $\mathcal{O}_K$ -basis of  $\mathcal{O}_L^\vee$ . We show that there exists a vector  $\mathbf{a} = (a_\tau)_{\tau \in \text{Gal}(L/K)} \in (\mathcal{O}_L^\vee)^\varphi$  such that, for any  $x \in L$ , we can write

$$f(x) = \sum_{\tau \in \text{Gal}(L/K)} a_\tau \cdot \tau(x).$$

To construct  $\mathbf{a}$ , we first construct another vector  $\mathbf{b} = (b_i)_{i \in [\varphi]} \in L^\varphi$  by setting  $b_i := f(e_i)$  for all  $i \in [\varphi]$ . We then define  $a_\tau := \mathbf{b}^T \cdot \tau(\mathbf{e}^\vee)$  for all  $\tau \in \text{Gal}(L/K)$ , where the automorphism  $\tau$  is applied component-wise. For any  $x = \sum_{i \in [\varphi]} x_i e_i \in L$  where  $x_i \in K$ , we observe that

$$\begin{aligned} \sum_{\tau \in \text{Gal}(L/K)} a_\tau \cdot \tau(x) &= \sum_{\tau \in \text{Gal}(L/K)} \mathbf{b}^T \tau(\mathbf{e}^\vee) \cdot \tau(x) \\ &= \sum_{\tau \in \text{Gal}(L/K)} \sum_{j \in [\varphi]} b_j \tau(e_j^\vee) \cdot \tau \left( \sum_{i \in [\varphi]} x_i e_i \right) \\ &= \sum_{\tau \in \text{Gal}(L/K)} \sum_{i,j \in [\varphi]} b_j \tau(e_j^\vee e_i) x_i \\ &= \sum_{i,j \in [\varphi]} b_j \text{Trace}_{L/K}(e_j^\vee e_i) x_i \\ &= \sum_{i \in [\varphi]} b_i x_i = \sum_{i \in [\varphi]} f(e_i) x_i = f(x). \quad \square \end{aligned}$$

**Lemma 12.** *Let  $L$  be a cyclotomic field,  $L/K$  a Galois extension,  $\mathfrak{f}_L$  be the conductor of  $L$ , and  $p$  be a rational prime with  $\mathfrak{f}_L < p$ . It holds that any  $\mathcal{O}_K/p\mathcal{O}_K$ -linear map  $f : \mathcal{O}_L/p\mathcal{O}_L \rightarrow \mathcal{O}_L/p\mathcal{O}_L$  can be expressed as an  $\mathcal{O}_L/p\mathcal{O}_L$ -linear combination of  $\text{Gal}(L/K)$ .*

*Proof.* In the proof of Lemma 11, we can see that any  $\mathcal{O}_K$ -linear map  $f : \mathcal{O}_L \rightarrow \mathcal{O}_L$  can be expressed as an  $\mathcal{O}_L^\vee$ -linear combination of  $\text{Gal}(L/K)$ . Since  $L$  is cyclotomic, it is known (see e.g. [LPR13]) that

$$\mathcal{O}_L \subseteq \mathcal{O}_L^\vee \subseteq \mathfrak{f}_L^{-1} \mathcal{O}_L.$$

Taking quotients, we have

$$\frac{\mathcal{O}_L}{p\mathcal{O}_L} \subseteq \frac{\mathcal{O}_L^\vee}{p\mathcal{O}_L} \subseteq \frac{\mathfrak{f}_L^{-1} \mathcal{O}_L}{p\mathcal{O}_L}.$$

However, since  $\mathfrak{f}_L < p$  and  $p$  is prime, we have

$$\frac{\mathcal{O}_L}{p\mathcal{O}_L} \subseteq \frac{\mathcal{O}_L^\vee}{p\mathcal{O}_L} \subseteq \frac{\mathfrak{f}_L^{-1} \mathcal{O}_L}{p\mathcal{O}_L} = \frac{\mathcal{O}_L}{p\mathcal{O}_L},$$

forcing  $\frac{\mathcal{O}_L^\vee}{p\mathcal{O}_L} = \frac{\mathcal{O}_L}{p\mathcal{O}_L}$ . The claim thus follows.  $\square$

## A.2 Vanishing Short Integer Solution Assumption

To prove the soundness of some of the argument systems proposed in this work, we rely on the vanishing short integer solution (vSIS) assumption, proposed in [CLM23] as a structured variant of the standard SIS assumption, and can be seen as a variant of the kRISIS assumption [ACL<sup>+</sup>22] without hints. To recall, the vSIS assumption is in fact a family of assumptions parametrised by a ring  $\mathcal{R}$ , a modulus  $q$ , a number of points  $n$ , a number of variables  $\mu$ , a norm bound  $\beta$  (for an implicit norm function), and a family  $\mathcal{G}$  of  $\mu$ -variate (possibly Laurent) polynomials over  $\mathcal{R}$ . It states that, given  $n$  randomly sampled evaluation points  $\mathbf{x}_i \leftarrow_{\$} (\mathcal{R}_q^\times)^\mu$  for  $i \in [n]$ , it is infeasible to find a polynomial  $g \in \mathcal{G}$  satisfying  $g(\mathbf{x}_i) = 0 \pmod q$  for all  $i \in [n]$  and whose coefficient vector has norm at most  $\beta$ .

Currently, the best known approach to solve a vSIS problem is to solve it as an unstructured SIS problem, except for extreme cases, e.g. when  $g$  is allowed to have a degree close to  $q$ . We refer to [CLM23] for more discussion about the conjectured hardness of vSIS.

Referring to the formulation of the vSIS assumption given in Definition 1, setting  $\chi$  to be the uniform distribution over  $\mathcal{R}_q^{n \times m}$  recovers the standard SIS assumption. More interestingly, by instantiating the distribution  $\chi$  differently, we can recover various variants of the vSIS assumption stated in the style of [CLM23]. For example, setting  $\chi$  to the uniform distribution of vectors of the form  $\bigotimes_{i \in [\mu]} (1, x_i)$  for  $x_i \in \mathcal{R}_q^\times$ , we recover the single-point multilinear variant. Another example is to set  $\chi$  to the uniform distribution of vectors of the form  $\bigotimes_{i \in [\mu]} (1, x^{2^i})$ , which corresponds to the single-point univariate variant.

## A.3 Reduction of Knowledge

In this paper we consider ternary relations  $\Xi \subseteq \{0, 1\}^* \times \{0, 1\}^* \times \{0, 1\}^*$ , where a tuple  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi$  consists of public parameters  $\text{pp}$ , statement  $\text{stmt}$  and witness  $\text{wit}$ . For presentation, we omit including  $\text{pp}$  when it is known from the context. We consider a modified and simplified definition of a reduction of knowledge [KP23] for the following reasons: All of our protocols are *public coin* and (*coordinate-wise*) *special sound* [FMN23] or similar.<sup>21</sup> Thus, public reducibility is automatic and we have (super-constant) sequential composition results due to known (tree) black-box extractors, whereas composition in [KP23] is limited a constant number of protocols. Lastly, we define a *relaxed* knowledge soundness notion which is not present in [KP23].

*Remark 12.* To turn soundness errors of probabilistic tests (such as Schwartz–Zippel) into knowledge errors, we merely need two uniformly random accepting transcripts. These are produced by (CW)SS extractors for example. We call such extractors 2-transcript extractors. We note that the protocols themselves are *not* (CW)SS, as not *any* pair of transcripts (with distinct challenges) suffices for extraction. However, given two transcripts (where the challenges are uniformly distributed conditioned on accepting), we can nonetheless bound the knowledge error. This occurs when extracting a reduction of knowledge whose soundness relies on a Schwartz–Zippel-type argument. All common extractors for (CW)SS satisfy the required distribution of transcripts, hence we can use these extractors as 2-transcript extractors. Importantly, we can “pretend” to deal with (CW)SS in terms of extracting two transcripts. Thus, the tree-special soundness is still applicable and the running time is bounded in the tree size (for state of the art extractors). Our definition of knowledge error is simply additive for sequential compositions of extractors. Hence we can compose as many extractors as we need, and the resulting extractor is efficient if the extracted tree of transcripts is remains polynomial in size.

**Definition 6 (Reduction of Knowledge (modified)).** *Let  $\Xi_0, \Xi_1$  be ternary relations. A reduction of knowledge  $\Pi$  from  $\Xi_0$  to  $\Xi_1$ , short  $\Pi: \Xi_0 \rightarrow \Xi_1$ , is defined by two PPT algorithms  $\Pi = (\mathcal{P}, \mathcal{V})$ , the prover  $\mathcal{P}$ , and the verifier  $\mathcal{V}$ , with the following interface:*

- $\mathcal{P}(\text{pp}, \text{stmt}_1, \text{wit}_1) \rightarrow (\text{stmt}_2, \text{wit}_2)$ : *Interactively reduce the input statement  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi_0$  to a new statement  $(\text{pp}, \widetilde{\text{stmt}}, \widetilde{\text{wit}}) \in \Xi_1$  or  $\perp$ .*
- $\mathcal{V}(\text{pp}, \text{stmt}) \rightarrow \widetilde{\text{stmt}}$ : *Interactively reduce the task of checking the input statement  $(\text{pp}, \text{stmt})$  w.r.t  $\Xi_0$  to checking a new statement  $(\text{pp}, \widetilde{\text{stmt}})$  w.r.t.  $\Xi_1$ .*

<sup>21</sup>See also Remark 12.

Let  $\langle \mathcal{P}, \mathcal{V} \rangle$  denote the interaction between  $\mathcal{P}$  and  $\mathcal{V}$ , as a function that takes as input  $(\text{pp}, \text{stmt}, \text{wit})$  and runs the prover  $\mathcal{P}$  (resp. verifier  $\mathcal{V}$ ) on input  $(\text{pp}, \text{stmt}, \text{wit})$  (resp.  $(\text{pp}, \text{stmt})$ ). At the end of the interaction,  $\langle \mathcal{P}, \mathcal{V} \rangle$  outputs the verifier's statement  $\widetilde{\text{stmt}}$  and prover's witness  $\widetilde{\text{wit}}$ . We define following properties.

**Definition 7 (Correctness).** *Let  $\Pi = (\mathcal{P}, \mathcal{V})$  be a reduction of knowledge from  $\Xi_0$  to  $\Xi_1$ . We say  $\Pi$  has correctness error  $\gamma(\cdot)$ , if for all  $(\text{pp}, \text{stmt}, \text{wit}) \in \Xi_0$*

$$\Pr[(\text{pp}, \widetilde{\text{stmt}}, \widetilde{\text{wit}}) \in \Xi_1 \mid (\widetilde{\text{stmt}}, \widetilde{\text{wit}}) \leftarrow \langle \mathcal{P}, \mathcal{V} \rangle(\text{pp}, \text{stmt}, \text{wit})] \geq 1 - \gamma(\text{pp}, \text{stmt}).$$

If  $\gamma \equiv 0$ , we call  $\Pi$  perfectly correct.

Our definitions of knowledge soundness and error are tailored to knowledge extractors for (coordinate-wise) special soundness (cf. Section A.4). Unlike [KP23], we require black-box extraction and ignore efficiency of the adversary.

**Definition 8 ((Black-Box) Knowledge Soundness).** *Let  $\Pi = (\mathcal{P}, \mathcal{V})$  be a reduction of knowledge. We say that  $\Pi$  is relaxed knowledge sound from  $\Xi_0^{KS}$  to  $\Xi_1^{KS}$  with knowledge error  $\kappa(\text{pp}, \text{stmt})$  if there exists a black-box expected polynomial-time extractor  $\mathcal{E}$  such that: For all  $\text{pp}, \text{stmt}$ , and every (unbounded) malicious prover  $\mathcal{P}^*$ , we have*

$$\begin{aligned} & \Pr[(\text{pp}, \text{stmt}, \text{wit}) \in \Xi_1^{KS} \mid \text{wit} \leftarrow \mathcal{E}^{\mathcal{P}^*}(\text{pp}, \text{stmt})] \\ & \geq \Pr[(\text{pp}, \widetilde{\text{stmt}}, \widetilde{\text{wit}}) \in \Xi_0^{KS} \mid (\widetilde{\text{stmt}}, \widetilde{\text{wit}}) \leftarrow \langle \mathcal{P}^*, \mathcal{V} \rangle(\text{pp}, \text{stmt})] - \kappa(\text{pp}, \text{stmt}). \end{aligned}$$

If  $\Pi$  is both correct and relaxed knowledge sound from  $\Xi_0^{KS}$  to  $\Xi_1^{KS}$ , then we say that  $\Pi$  is knowledge sound from  $\Xi_0^{KS}$  to  $\Xi_1^{KS}$ .

We assert no composition theorem. Instead we appeal to the fact that all of our protocols are (coordinate-wise) special sound, so we eventually require are tree-CWSS extractor. These are known to exist, even for the Fiat–Shamir transformed protocol, see e.g. [FMN23]. Formally, we must translate reductions of knowledge to proofs of knowledge to apply special soundness. But this is a triviality:

**Definition 9.** *Let  $\Pi = (\mathcal{P}, \mathcal{V})$  be a reduction of knowledge from  $\Xi_0$  to  $\Xi_1$ . We define the induced proof of knowledge  $\widehat{\Pi} = (\widehat{\mathcal{P}}, \widehat{\mathcal{V}})$  of  $\Pi$ , where the prover  $\widehat{\mathcal{P}}$  sends an additional (final) protocol message  $\widehat{\text{wit}}$ , and the verifier outputs the bit  $(\text{pp}, \widetilde{\text{stmt}}, \widehat{\text{wit}}) \in \Xi_1$ .*

By considering the induced proof of knowledge for  $\Pi: \Xi_0 \rightarrow \Xi_1$ , our for  $\Pi_{KS}: \Xi_0^{KS} \rightarrow \Xi_1^{KS}$  in case of relaxed knowledge soundness, we see that the notion of correctness and (relaxed) knowledge soundness is equivalent to the respective notion for reduction of knowledge, with the same knowledge error. Hence, we can indeed apply all our tools for (coordinate-wise special sound) proofs of knowledge.

#### A.4 Coordinate-Wise Special Soundness

We recall the notion of coordinate-wise special soundness (CWSS) from [FMN23] in a simple form. Let  $S$  be a finite set and  $\ell \geq 1$ . For  $i \in [\ell]$  define the following relation  $\equiv_i$  for two vectors  $\mathbf{x}, \mathbf{y} \in S^\ell$ :

$$\mathbf{x} \equiv_i \mathbf{y} \iff x_i \neq y_i \quad \text{and} \quad x_j = y_j \quad \forall j \in [\ell] \setminus \{i\}.$$

This means that  $\mathbf{x}$  and  $\mathbf{y}$  differ in exactly the  $i$ -th coordinate. Next, we consider the following set:

$$\Gamma(S, \ell) := \{(\mathbf{x}_0, \dots, \mathbf{x}_\ell) \subseteq (S^\ell)^{\ell+1} : \forall i \in [\ell], \mathbf{x}_0 \equiv_i \mathbf{x}_i.\}$$

In other words,  $(\mathbf{x}_0, \dots, \mathbf{x}_\ell) \in \Gamma(S, \ell)$  if for every coordinate  $i \in [\ell]$ , there exists exactly one vector  $\mathbf{x}_i$  that differs from  $\mathbf{x}_0$  in exactly (and only) the  $i$ -th coordinate. One can think of  $\mathbf{x}_0$  as the ‘‘central’’ vector.

Intuitively, coordinate-wise special soundness for three-round interactive proofs says that given  $\ell + 1$  valid transcripts with challenges  $\mathbf{x}_0, \dots, \mathbf{x}_\ell$  which satisfy  $(\mathbf{x}_0, \dots, \mathbf{x}_\ell) \in \Gamma(S, \ell)$ , one can efficiently extract the witness. In this paper, we will call such protocols  $\ell$ -CWSS. To argue knowledge soundness from coordinate-wise special soundness in the context of reduction of knowledge, we will use the following lemma from [FMN23].

**Lemma 13 (CWSS).** *Let  $\ell \in \mathbb{N}$ , and  $S$  be a finite set of cardinality  $N$ . Let  $\mathcal{C} := S^\ell$  and take a verification function  $\mathcal{V} : \mathcal{C} \times \{0, 1\}^* \rightarrow \{0, 1\}$ . Then there exists an extractor algorithm  $\mathcal{E}$ , which given oracle access to a probabilistic algorithm  $\mathcal{A}$  such that*

$$\varepsilon := \Pr[\mathcal{V}(\mathbf{x}, \mathcal{A}(\mathbf{x})) = 1],$$

*where the probability is over the choice of  $\mathbf{x} \leftarrow \mathcal{C}$  and random coins of  $\mathcal{A}$ , it makes an expected number of at most  $\ell + 1$  queries to  $\mathcal{A}$  and with probability at least*

$$\varepsilon - \ell/N$$

*outputs  $\ell + 1$  pairs  $(\mathbf{x}_i, y_i)_{0 \leq i \leq \ell}$  such that  $V(\mathbf{x}_i, y_i) = 1$  for all  $0 \leq i \leq \ell$  and  $\{\mathbf{x}_0, \dots, \mathbf{x}_\ell\} \in \Gamma(S, \ell)$ .*