

# Concretely Efficient Zero-Knowledge Argument System Based on Input Transformation and Direct Computation

Frank Y.C. Lu

YinYao Inc.

**Abstract.** We introduce a new transparent zero-knowledge argument system based on the direct computation concept described in this paper. Our protocol converts input parameters into a format that any circuit can process directly. Once the circuit output can be computed using transformed inputs, the output of the polynomial a circuit generates can also be correctly computed by the verifier, allowing us to reduce the polynomial evaluation cost significantly. In the default setting, the prover runtime cost for group exponentiation operations is only the square root of the degree ( $\sqrt{p_d}$ ) of the polynomial the circuit generates. Furthermore, leveraging the “merging through addition” and “bootstrapping with breakers” techniques, the size of the polynomial our protocol generates can be much smaller than the total number of multiplicative operations. Our benchmark result shows our approach can significantly improve both prover runtime and verifier runtime performance over state-of-the-art by almost or over one order of magnitude while keeping the communication cost comparable with that of the state-of-the-art. Our approach also allows our protocol to be memory-efficient while being non-interactive. The theoretical memory cost of our protocol is  $O(b)$ , where  $b = \sqrt{p_d}$  in the default setting. Lowering the  $b$  value will result in better prover runtime performance at the expense of the higher communication cost.

## 1 Introduction

Ever since the discoveries of interactive proofs (IPs) [26] and probabilistically checkable proofs (PCPs) [6] [5] [4] [3] in the late last century, there has been a tremendous amount of research in the area of proof systems. More recently, the rise of blockchain and Web3 has finally triggered real-world interest in zero-knowledge systems.

Due to the expensive computation cost in the setup phase of earlier SNARKs (Succinct Non-Interactive Argument of Knowledge), the industry developed protocols that have the structured reference string (SRS) be constructible in a “universal and updatable” fashion. The first such universal SNARK was in Groth et al. [27], and Maller et al. improved the SRS size from quadratic to linear in Sonic [30]. More recently developed protocols such as PLONK [24], MARLIN [17] are universal fully-succinct SNARK with significantly improved prover runtime compared to the fully-succinct Sonic. However, many of these universal succinct SNARKs systems require trusted setup, and the prover run-time of these protocols is prohibitively expensive even with the latest improvements such as HyperPlonk [16], usually takes over 100 seconds on a single-threaded CPU for a circuit with over  $2^{20}$  constraints.

Other classes of protocols including the Goldwasser, Kalai, and Rothblum (GKR) class such as Hyrax [35], Virgo [40]; MPC-in-the-head class of Kushilevitz, Ostrovsky, and Sahai such as ZKBoo [25] and Liger/Liger++ [1] [9] offer efficient prover runtimes that are at least one order of magnitude more efficient than pairing-based SNARKs. However, these protocols are largely ignored by the industry (e.g., the blockchain community) due to their expensive verifier runtime and high communication cost (hundreds of KBs) compared to fully succinct SNARK and STARK [8] protocols.

NIZKs such as SpartanNIZK [33] and later Lakonia [34] seem to offer a much more balanced approach, where they offer efficient prover runtime (6-18 seconds single thread) and competitive communication costs for large circuits while not being layer dependent. However, the downside of these

protocols is that their verifier performance is still expensive, usually in the 400+ ms range on a single-threaded CPU.

Recent advances in multivariate polynomial commitment schemes (PCS) such as Binius [20] can, in theory, improve the prover runtime performance for protocols using multivariate PCS even better with the right hardware. However, it does not bridge the gap in communication cost with popular ECC-based PCS.

Another category of protocols emphasizes on memory-efficiency such as garbled circuits [29] [23] [28] and Vector Oblivious Linear Evaluation (VOLE) protocols [12] [32] [14] [13] [39] [36] [7] [38] generally offer better prover performance. However, their verifier runtimes are just as expensive and generally require a designated verifier with very expensive communication cost.

Our aim is to create a new transparent zero-knowledge protocol that is easy to code for both developers and business people and, at the same time, offers the best overall performance while free of significant performance shortcoming in any area. Specifically, we want to keep the communication cost comparable to those of the state-of-the-art and greatly improve both the prover runtime and the verifier runtime costs. Finally, we also want our new system to be memory-efficient without requiring a designated verifier like that of VOLE protocols.

## 2 Summary of Contributions

Our approach is to design a new class of protocols that allows verifiers to validate circuit outputs by directly examining circuit inputs without going through some intermediate translation phase. In our protocol, circuit inputs in the Pedersen commitment form are converted to linear polynomials in  $\mathbb{F}_q$  so that verifiers can use standard arithmetic operations in  $\mathbb{Z}_q$  to just “execute” the evaluating circuit to get its output. Our protocol does not require trusted setup and depends only on discrete logarithm assumption.

There were past attempts that somewhat enabled verifiers to “validate” each multiplication gate on its own, such as Cramer and Damgård [19] and more recent designated-verifier (which is a limitation itself) VOLE (LPZK in particular) [22] [38] [21] [37] based protocols. In these older strategies, each multiplication gate computation is actually not computed but “confirmed” by the verifier using transcripts tied to each multiplication gate. As a result, the communication/verifier costs of these earlier protocols are generally linear in the number of multiplication gates in a circuit.

On the other hand, the input transformation technique introduced by our protocol allows verifiers to use transformed inputs to directly (one operation for each operation, like we do with clear text data) compute the circuit (important), and the verifier computed output is still bound to the challenge  $x$ . This is a first and brings us three direct benefits; 1) After “computing” the whole circuit using transformed inputs, the verifier can now validate sub-linear sized proof transcripts in sub-linear runtime. 2) Since the whole circuit is linearly/directly computed by the verifier, we can break a large circuit into several smaller sub-circuits, minimizing the memory footprint to that of a sub-circuit. 3) Because the circuit is linearly “computed” by both the prover and the verifier, it gives the developer the power to significantly reduce the size of the circuit by combining “other” protocols in the middle and bypassing the “inactive” part of the circuit logic.

In our protocol, committed input parameters in  $\mathbb{G}$  are transformed to linear polynomials in  $\mathbb{Z}_q$ . For simple, circuit  $a_1^d + a_2^d + a_3^d = r$  takes inputs  $a_1, a_2, a_3$  and outputs  $r$ . In our protocol, inputs  $a_1, a_2, a_3$  and output  $r$  are committed by the prover using Pedersen commitment. The prover also provides the transformed inputs  $a_1, a_2, a_3$  in linear polynomial form  $a'_1, a'_2, a'_3 \in \mathbb{Z}_q$  s.t.  $a'_i = a_i + X\alpha_i \in \mathbb{Z}_q$  ( $\alpha_i$  is its blinding key and  $x$  is the challenge generated during runtime). Since the transformed inputs are in  $\mathbb{F}_q$ , the verifier can plug these values directly into a circuit and just “execute” them to get an output  $o$  e.g.  $a_1'^d + a_2'^d + a_3'^d = o$ . The circuit output  $o \in \mathbb{Z}_p$  is the evaluation at point  $x$  of a degree  $d$  polynomial s.t.  $f(x) = o$ . The constant term of this polynomial is the circuit output  $r$  and all other coefficients are blinding keys. If the prover can prove 1) it knows all coefficients of the output polynomial before

the evaluation point is given (e.g., using a polynomial commitment) and 2) all input transformations are legit, then we say the proof is valid.

The output polynomial in the example above has degree  $d$  because the transformed inputs (linear polynomials) have degree 1. Taking them to their  $d$ th power and adding them together will get a polynomial of degree  $d$ . So if the circuit is something like  $a_1^3 + a_2^3 + a_3^3 + \dots + a_t^3 = o$ , the output polynomial has a degree of 3 regardless of what the total number of addition and multiplication operations in that circuit is. Throughout our paper, we use the symbol  $p_d$  (short for “polynomial degree”) to denote the degree of the final polynomial that leads to the circuit output.  $p_d$  is different from “multiplication depth” or the “total number of multiplications in a circuit”. For example, for a circuit  $a_1^3 \cdot a_2^5 + a_3^6 = r$ ,  $s = 7$  ( $a_1^3 \cdot a_2^5$  is a polynomial of degree 8), which is bigger than the multiplication depth (5) but smaller than the total number of multiplications (12).

A key advantage of our protocol is that if the transformed inputs  $a'_i = a_i + x\alpha_i \in \mathbb{Z}_q$  are known and can be validated, the verifier can just execute the circuit to get the correct polynomial output  $f(x) = o$  for some evaluation point  $x$ , allowing us to bypass the need to use an expensive polynomial commitment scheme to evaluate the whole polynomial.

**High performance** Unlike other zero-knowledge protocols that depend on polynomial commitment, the result of evaluation point  $x$  in our protocol is computed by the verifier (through direct computation of transformed inputs in  $\mathbb{Z}_q$ ) itself, so it has to be accurate. In addition, the constant term coefficients are committed by the verifier. This allows us to bypass the expensive PCS protocols and only evaluate a small polynomial of  $O(b)$  degrees, where  $b$  is a parameter set by the user that tells the protocol where to “reset” a degree  $d$  polynomial back to a linear polynomial (degree 1), a technique used to slow the growth of polynomial degree. In the default setting, we set  $b = \sqrt{p_d}$ .

Specifically, when processing a high-depth circuit of  $2^{20}$  sequential multiplication gates ( $s = n$ ) with 20 inputs on a single CPU thread, the performance of our protocol is: 0.57 seconds for the prover runtime cost; 12 milliseconds for the verifier runtime cost; and 30 kilobytes for the communication cost. To the best of our knowledge, our protocol offers the best prover/verifier runtime performance in the literature by a large margin.

In the memory-efficient setting, the theoretical memory cost of our protocol is  $O(b)$ . This makes our protocol extremely attractive because VOLE-based memory-efficient protocols generally require one round of interaction and are extremely expensive in terms of verifier runtime cost and communication cost.

**Same Format for Circuit Inputs and Outputs** Having both inputs and output(s) in the same format (Pedersen Commitment) allows the output(s) of one circuit to be directly reused as inputs to another circuit. This is a really useful feature in practice when verifying data in a publicly accessible/verifiable data store (not limited to blockchain). e.g., allows many participants to continuously manage/update a datastore as long as they can prove these updates were correctly computed from existing data. While other zero-knowledge protocols may be able to support such a feature in theory by mapping witnesses to some publicly accessible committed or encrypted data, it does not come naturally and will require additional cost that is not accounted for.

Since both inputs and outputs are in the same formats: Pedersen commitment in  $\mathbb{G}$  (standard format) and linear polynomial in  $\mathbb{F}_q$  (transformed format). Since data are perfectly hidden under both formats, results of direct computation are automatically perfectly hidden. Therefore, our protocol is zero knowledge as long as we can ensure transcripts used for transformations between two formats are zero knowledge.

**Ease of Use** Another less appreciated benefit of the input transformation based approach is that developers can just wire a circuit in a much more straight forward fashion than the expensive arithmetization process used in other protocols.

We introduce our protocol in an interactive setting where all verifier challenges are random field elements. In practice, we assume the Fiat-Shamir heuristic is applied to our protocol to obtain a non-interactive zero-knowledge argument in the random oracle model.

### 3 Preliminaries

#### 3.1 Assumption

**Definition 1.** (Discrete Logarithmic Relation) For all PPT adversaries  $\mathcal{A}$  and for all  $n \geq 2$  there exists a negligible function  $\mathit{negl}(\lambda)$  s.t.

$$\Pr \left[ \begin{array}{l} \mathbb{G} = \mathit{Setup}(1^\lambda), g_0, \dots, g_{n-1} \xleftarrow{\$} \mathbb{G} \\ a_0, \dots, a_{n-1} \in \mathbb{Z}_p \\ \leftarrow \mathcal{A}(\mathbb{G}, g_0, \dots, g_{n-1}) \end{array} \middle| \begin{array}{l} \exists a_i \neq 0 \\ \wedge \prod_{i=0}^{n-1} g_i^{a_i} = 1 \end{array} \right] \leq \mathit{negl}(\lambda)$$

The Discrete Logarithmic Relation assumption states that an adversary can't find a non-trivial relation between the randomly chosen group elements  $g_0, \dots, g_{n-1} \in \mathbb{G}^n$ , and that  $\prod_{i=0}^{n-1} g_i^{a_i} = 1$  is a non-trivial discrete log relation among  $g_0, \dots, g_{n-1}$ . Please note the generators we use in this paper are  $g, h, u \in \mathbb{G}$ .

#### 3.2 Zero-Knowledge Argument of Knowledge

Interactive arguments are interactive proofs in which security holds only against computationally bounded provers. In an interactive argument of knowledge for a relation  $\mathcal{R}$ , a prover convinces a verifier that it knows a witness  $w$  for a statement  $x$  s.t.  $(x, w) \in \mathcal{R}$  without revealing the witness itself to the verifier.

Let  $(\mathcal{P}, \mathcal{V})$  denote a pair of PPT interactive algorithms, and **Setup** denotes a non-interactive setup algorithm that outputs public parameters  $pp$  given a security parameter  $\lambda$ . Let  $\langle \mathcal{P}(pp, x, w), \mathcal{V}(pp, x) \rangle$  denote the output of  $\mathcal{V}$  on input  $x$  after its interaction with  $\mathcal{P}$ , who has knowledge of witness  $w$ . The triple  $(\mathit{Setup}, \mathcal{P}, \mathcal{V})$  is called an argument for relation  $\mathcal{R}$  if for all non-uniform PPT adversaries  $\mathcal{A}$  it satisfies completeness, soundness, and zero-knowledge definitions defined below:

**Definition 2.** (Perfect Completeness) The triple  $(\mathit{Setup}, \mathcal{P}, \mathcal{V})$  satisfies perfect completeness if for all PPT  $\mathcal{A}$ :

$$\Pr \left[ \begin{array}{l} (pp, x, w) \notin \mathcal{R} \text{ or} \\ \langle \mathcal{P}(pp, x, w), \mathcal{V}(pp, x) \rangle = 1 \end{array} \middle| \begin{array}{l} pp \leftarrow \mathit{Setup}(1^\lambda) \\ (x, w) \leftarrow \mathcal{A}(pp) \end{array} \right] = 1$$

The soundness notion we consider in this work is computational witness-extended emulation.

**Definition 3.** (Computational Witness-Extended Emulation or CWEE) Given a public-coin interactive argument tuple  $(\mathit{Setup}, \mathcal{P}, \mathcal{V})$  and arbitrary prover algorithm  $\mathcal{P}^*$ , let **Recorder**  $(\mathcal{P}^*, pp, x, s)$  denote the message transcript between  $\mathcal{P}^*$  and  $\mathcal{V}$  on shared input  $x$ , initial prover state  $s$ , and  $pp$  generated by **Setup**. Furthermore, let  $\mathcal{E}$  **Recorder**  $(\mathcal{P}^*, pp, x, s)$  denote a machine  $\mathcal{E}$  with a transcript oracle for this interaction that can rewind to any round and run again with fresh verifier randomness. The tuple  $(\mathit{Setup}, \mathcal{P}, \mathcal{V})$  has CWEE if for every deterministic polynomial time  $\mathcal{P}^*$  there exists an expected polynomial time emulator  $\mathcal{E}$  s.t. for all non-uniform polynomial time adversaries  $\mathcal{A}$  the following holds:

$$\Pr \left[ \mathcal{A}(tr) = 1 \middle| \begin{array}{l} pp \leftarrow \mathit{Setup}(1^\lambda) \\ (x, s) \leftarrow \mathcal{A}(pp) \\ tr \leftarrow \mathbf{Recorder}(\mathcal{P}^*, pp, x, s) \end{array} \right] -$$

$$Pr \left[ \begin{array}{l} \mathcal{A}(tr) = 1 \wedge \\ tr \text{ accepting} \\ \implies (x, w) \in \mathcal{R} \end{array} \middle| \begin{array}{l} pp \leftarrow Setup(1^\lambda) \\ (x, s) \leftarrow \mathcal{A}(pp) \\ (tr, w) \leftarrow \mathcal{E}^{Recorder(\mathcal{P}^*, pp, x, s)}(pp, x) \end{array} \right] \leq \mathit{negl}(\lambda)$$

Informally, if an adversary can produce an argument that satisfies the verifier with some probability, then there exists an emulator producing an identically distributed argument with the same probability, as well as a witness. The zero-knowledge property requires that the verifier doesn't learn anything about the witness from its interaction with an honest prover.

**Definition 4.** (Perfect Special Honest Verifier Zero Knowledge for Interactive Arguments) An interactive proof is  $(Setup, \mathcal{P}, \mathcal{V})$  is a perfect special honest verifier zero knowledge (PHVZK) argument of knowledge for  $\mathcal{R}$  if there exists a probabilistic polynomial time simulator  $\mathcal{S}$  such that all interactive adversaries  $\mathcal{A}$  have the following property for every  $(x, w, \sigma) \leftarrow \mathcal{A}(pp) \wedge (pp, x, w) \in \mathcal{R}$ , where  $\sigma$  stands for verifier's public coin randomness for challenges

$$Pr \left[ \begin{array}{l} \mathcal{A}(tr) = 1 \\ \left| \begin{array}{l} pp \leftarrow Setup(1^\lambda), \\ tr \leftarrow \langle \mathcal{P}(pp, x, w), \mathcal{V}(pp, x) \rangle \end{array} \right. \end{array} \right] = \\ Pr \left[ \begin{array}{l} \mathcal{A}(tr) = 1 \\ \left| \begin{array}{l} pp \leftarrow Setup(1^\lambda), \\ tr \leftarrow \mathcal{S}(pp, x, \sigma) \end{array} \right. \end{array} \right]$$

Above property states that the adversary chooses a distribution over statements  $x$  and witnesses  $w$  but is not able to distinguish between the simulated transcripts and the honestly generated transcripts for a valid statement/witnesses pair, and that the simulator has access to the randomness used by the verifier.

**Definition 5.** (Public Coin) All messages sent from  $\mathcal{V}$  to  $\mathcal{P}$  are chosen uniformly at random and independently of  $\mathcal{P}$ 's messages.

### 3.3 Zero Knowledge Proof of Discrete Logarithm

For a prover to prove it has the knowledge of a discrete logarithmic  $\kappa$  of some group element  $s = g^\kappa \in \mathbb{G}$ . We define the relation for this protocol as  $\mathcal{R}_{PoD} = \{(h, s; \kappa) : s = g^\kappa\}$ . We also define two functions ( $ProveDL, VerifyDL$ ) for provers and verifiers to create and verify proof transcripts:

- $ProveDL(g, \kappa) \rightarrow tr_\kappa$  generates the proof transcript  $tr_\kappa$ , where  $\kappa$  is the witness.
- $VerifyDL(g, s, tr_\kappa) \rightarrow b \in \{0, 1\}$  takes a proof transcript  $tr_\kappa$  and a pair of group elements with discrete log relation  $(g, s \in \mathbb{G} \wedge s = h^\kappa)$ , and outputs *true* if the knowledge of the relation is verified, *false* otherwise.

In this paper, we assume the underlying implementation of the proof of discrete logarithm protocol is Schnorr's protocol [31]. We know for a fact that Schnorr's protocol has perfect completeness, special honest verifier zero knowledge, and computational witness-extended emulation.

### 3.4 Notations

Let  $\mathbb{G}$  denote any type of secure cyclic group of prime order  $p$ , and let  $\mathbb{Z}_p$  denote an integer field modulo  $p$ . Group elements other than generators are denoted by capital letters. e.g.,  $C = u_1^{a_1} u_2^{a_2} \dots u_n^{a_n} \in \mathbb{G}$  is a commitment committed to a vector  $\vec{a}$  denoted by a capital letter, and  $B \in \mathbb{G}$  is a random group element also denoted by a capital letter. For generators used as base points to compute other group elements in our protocol, such as  $\vec{g}, h \in \mathbb{G}$ , we use lower case letters to denote them. Greek letters are used to label hidden key values. e.g.,  $v$  is the blinding key for Pedersen commitment  $P$  on generator  $h \in \mathbb{G}$  s.t.  $P = g^a h^v$ . Finally, we use standard vector notation  $\vec{v}$  to denote vectors. i.e.,  $\vec{a} \in \mathbb{Z}_p^n$  is a list of  $n$  values  $a_i$  for  $i = \{1, 2, \dots, n\}$  in  $\mathbb{Z}_p$ .

We write  $\mathcal{R} = \{(Public\ Inputs ; Witnesses) : Relation\}$  to denote the relation  $\mathcal{R}$  using the specified public inputs and witnesses.

## 4 Protocol for Arbitrary Circuits

We first define the relation for the base version of our protocol. For  $l$  input parameters, let  $\mathcal{C}_{\mathbb{F}}$  represent the set of arbitrary arithmetic circuits in  $\mathbb{F}$ , there exists a zero knowledge argument for the relation:

$$\begin{aligned} \{(g, h, u, \vec{P}, R \in \mathbb{G}, E \in \mathcal{C}_{\mathbb{F}}; \vec{a}, \vec{v}, r, \phi \in \mathbb{Z}_p) : E(\vec{a}) = r \\ \wedge P_i = g^{a_i} h^{v_i} \wedge R = g^r h^{\phi}\} \end{aligned} \quad (1)$$

$g, h, u$  are initial public parameters  $pp$  generated during setup. The above relation states that each input parameter to a circuit is represented by a commitment  $P_i$  in  $\mathbb{G}$ , which hides each input value  $a_i$  with a blinding key  $v_i$ .  $r$  is the output of circuit  $E$  computed from inputs  $\vec{a}$ , which is also committed in  $R \in \mathbb{G}$  with blinding key  $\phi$ .

The main idea behind the “input transformation” concept is the process of transforming committed inputs in  $\mathbb{G}$  to linear polynomials in  $\mathbb{F}$  so that the verifier can perform addition and multiplication operations “as is”. For an input commitment  $P_i$  s.t.  $P_i = g^{a_i} h^{v_i} \in \mathbb{G}$  where  $a_i$  is the input value and  $v_i$  is its blinding key, its transformed value linear polynomial is  $a_i' \in \mathbb{Z}_q$  is :

$$a_i' = a_i + x \cdot \alpha_i \in \mathbb{Z}_q \quad (2)$$

$x$  is the challenge provided by the verifier during runtime, and the blinding key of each input is replaced by a random  $\alpha_i$  s.t.  $\alpha_i \neq v_i$ . Likewise, the circuit output commitment  $R = g^r h^{\phi} \in \mathbb{G}$  also has a matching linear polynomial in  $\mathbb{Z}_q$  with blinding key  $\epsilon$ .

$$r' = r + x \cdot \epsilon \in \mathbb{Z}_q \quad (3)$$

Where  $r'$  is “directly computed” from a list of inputs  $a_i'$  ( $i = \{1, \dots, l\}$ ) by the verifier. Since inputs are in  $\mathbb{Z}_q$ , verifiers can perform arithmetic operations on them just as they do on decrypted numbers. The output value of a circuit evaluation is now a polynomial with  $p_d + 1$  degrees evaluated at point  $x$ . The constant term of the output polynomial is the circuit output  $r$ , and the coefficient of the degree one term is the blinding key  $\epsilon$  of the circuit output.

In the next two sub-sections, we explain our protocol in two steps: In the first step, we introduce a sub-protocol (Protocol InputMapping) that allows the prover to prove each input in  $\mathbb{G}$  is correctly transformed to that in  $\mathbb{Z}_q$ ; In the second step, we introduce the full protocol (Protocol AriCircuit) that uses the aforementioned sub-protocol to validate transformed inputs and proves the circuit output is correctly computed from circuit inputs as relation 15 states. Before jumping into details, there are three key techniques that our protocol leverages that are important to go over:

**Use of two primes** We use two prime domains  $p, q$  in our protocol. All transformed inputs are in field  $\mathbb{Z}_q$ , where the verifier can perform circuit-related computations, and commitments are in  $\mathbb{G}$  of order  $p$ . There are two simple reasons for using  $\mathbb{Z}_q$  in circuit computation: 1) Prime  $p$  is not NTT friendly. 2) It is cheaper to work in a smaller field (we default  $q$  to a 61-bit prime).

**Hiding data in  $\mathbb{Z}_q$  with a significantly larger value in  $\mathbb{F}$**  We use significantly larger numbers in  $\mathbb{F}$  to hide data in  $\mathbb{Z}_q$ . For example, if we want to hide some 61-bit value  $r \in \mathbb{Z}_q$ . The prover can use a 141-bit random value  $\omega$  to hide it s.t.  $r + \omega \in \mathbb{F}$  will perfectly hide  $r$  except for a negligible chance of at least  $2^{-80}$  (e.g., when  $\omega$  is too large or too small, the prover just picks another  $\omega$ ).

The reason we use this hack is because we often perform operations on committed values in  $\mathbb{G}$  with order  $p$ . If the committed value overflows  $p$  after some operation (e.g., for  $r, x \in \mathbb{Z}_p$ , the exponent of  $g^{rx}$  will almost certain overflow  $p$ ), then we cannot use mod  $q$  on the opening of the commitment to get  $rx \in \mathbb{Z}_q$  without changing  $r$  to something else.

**Polynomial degree of a circuit** As mentioned earlier, the degree size of the polynomial generated by our protocol is not necessarily tied to the number of multiplications. In most real-world use cases, we believe the  $p_d$  value generated by our protocol is smaller than the total number of multiplications.

This is because every time an “addition or subtraction” operation is performed, the  $p_d$  value of two input wires merges. For example, if the first input wire  $a_1^7$  has  $p_d = 7$  and the second input wire  $a_1^6$  has  $p_d = 6$ , their sum is an output wire with merged  $p_d = 7$ . The more addition/subtraction operations we have in a circuit, the smaller the  $p_d$  value compared to gate count. The sum-of-products circuit explained earlier only has a multiplication depth of 1.

Another technique we use to keep the  $p_d$  value small is using a technique called “breakers” to bootstrapping the  $p_d$  value by breaking a circuit into multiple smaller sub-circuits, a concept we will detail in Section 4.2.

#### 4.1 The Sub-Protocol for Linear Polynomial to Pedersen Commitment Mapping Validation

We define a sub-protocol that validates committed inputs in  $\mathbb{G}$  is correctly mapped to transformed inputs in  $\mathbb{Z}_q$ , which is defined by the following relation:

$$\{(g, h, \vec{P} \in \mathbb{G}, \vec{a}' \in \mathbb{Z}_q; \vec{a}, \vec{\alpha} \in \mathbb{Z}_q, \vec{v} \in \mathbb{Z}_p) : P_i = g^{a_i} h^{v_i} \wedge a'_i = a_i + X\alpha_i\} \quad (4)$$

$X$  is the challenge generated during runtime, so  $\vec{a}'$  is only available during runtime. The relation above says that for any commitment  $P_i$  of value  $a_i$  and blinding key  $v_i$ , the prover will provide  $a'_i$  during runtime that s.t.  $a'_i = a_i + x\alpha_i$  for some random public coin challenge  $x$ .

In this relation, we assume the committed inputs are in Pedersen commitment format because it has simple representation and wide adaptability. However, this is not required. For example, one can develop a commitment-to-linear polynomial mapping validation algorithm for a post-quantum commitment.

We make sure the domain of transformed input in  $\mathbb{Z}_q$  is friendly to NTT. When multiplying two polynomials of degree  $p_d$ , the trivial approach to compute the output polynomial’s coefficients would require a runtime cost of  $O(p_d^2)$ , whereas the NTT-based approach would reduce that to  $O(p_d \log p_d)$ . Prime  $q$  is also expected to be a lot shorter than  $p$  in bits, therefore improving the runtime cost while at the same time lowering the communication cost.

**Setup Phase** Before the random challenge  $x$  is available, the prover commits to two of the group elements  $S_i, T_i$  for each  $i$ th Pedersen input.

The first commitment  $S_i$  commits to the randomly generated blinding keys  $\omega_i \in \mathbb{Z}_p$  used in transformation. Note  $S_i$ ’s blinding key  $v_i$  is also the input commitment  $\vec{P}_i$  but on a different generator, and the committed witness is  $w \cdot q$  instead of  $w$ , we will see shortly why this is needed.

$$S_i = g^{\omega_i \cdot q} u^{v_i} \in \mathbb{G} \quad i = \{1, \dots, l\} \quad (5)$$

The second commitment  $T_i$  commits to the new blinding key  $\alpha_i$  used to hide data in transformed input  $a'_i \in \mathbb{Z}_q$ . We do that by committing to the difference between each  $v_i$  and  $\alpha_i$ , which validates the mapping between the Pedersen input’s blinding key  $v_i$  and that of the transformed input in  $\mathbb{Z}_q$ .

$$T_i = g^{v_i - \alpha_i} \in \mathbb{G} \quad i = \{1, \dots, l\} \quad (6)$$

The setup phase of the protocol is detailed below. This part is called before the random challenge  $x$  is generated.

<i>Input</i> : $(; \vec{a}, \vec{\alpha}, \vec{\omega} \in \mathbb{Z}_q, \vec{v} \in \mathbb{Z}_p)$	
<i>P's input</i> : $(; \vec{a}, \vec{\alpha}, \vec{\omega}, \vec{v});$	
<i>P compute</i> :	
$S_i = g^{\omega_i \cdot q} u^{v_i} \in \mathbb{G}$	$i = \{1, \dots, l\}$
$T_i = g^{v_i - \alpha_i} \in \mathbb{G}$	$i = \{1, \dots, l\}$
$\mathcal{P} \rightarrow \mathcal{V} : \vec{S}, \vec{T}$	

Protocol InputMapping - Setup

Once the setup phase completes, the prover then sends  $S_i, T_i$  for  $i = \{1, \dots, l\}$  to the verifier.

**Validation Phase** After the random challenge  $x$  is generated, the prover computes  $\vec{a}'$  and sends them to the verifier. Next, the protocol checks the mapping between transformed inputs in  $\mathbb{Z}_q$  to those in group  $\mathbb{G}$ .

For each  $a'_i$ , the prover provides transcript  $e_i$ , which is used to convert the blinding element  $x\alpha \in \mathbb{Z}_q$  of each  $a'_i$  to its raw form  $x\alpha \in \mathbb{F}$  (without mod  $q$ ).

$$e_i = ((x\alpha_i \bmod q) - x\alpha_i) \cdot x + \omega_i \cdot q \in \mathbb{F} \quad i = \{1, \dots, l\} \quad (7)$$

The part of  $e_i$  on the left of the addition sign  $((x\alpha_i \bmod q) - x\alpha_i) \cdot x$  is around  $\approx 183$ -bits, which is small enough s.t. an adversary can use brutal force attack to extract  $\alpha_i$ . We can rewrite this left part to  $s_i \cdot q$  for some  $|s_i| \leq 122$ -bits and equation 7 to  $e_i = (s_i + \omega_i) \cdot q$ . The prover can use a 183-bit random value  $\omega_i$  (80-bits larger than  $s_i$ ) to perfectly hide it except for a negligible probability of at most  $2^{-80}$ . Note this makes  $|e_i| \leq |\omega| + |q| \approx 244$ -bits.

With transcripts  $S_i, T_i, e_i$ , the verifier can compute a new commitment that closely matches the structure of  $P_i$ .

$$g^{a'_i \cdot x - e_i} \cdot T_i^{x^2} \cdot S_i = (g^x)^{a_i + xv_i} u^{v_i} \in \mathbb{G} \quad i = \{1, \dots, l\} \quad (8)$$

With  $(g^x)^{a_i + xv_i}$  available, we can compute a “public key”  $PK_{v_i}$  for witness  $v_i$ :

$$(h^x / (g^{x^2} u))^{v_i} = \frac{P_i^x}{g^{a'_i \cdot x - e_i} \cdot T_i^{x^2} \cdot S_i} \in \mathbb{G} \quad i = \{1, \dots, l\} \quad (9)$$

Where  $(h^x / (g^{x^2} u))$  is the generator for witness  $v_i$ . Using another challenge  $k \in \mathbb{Z}_p$ , we can create a public key for validating all  $P_i$  to  $a'_i$  mappings at once:

$$PK_{v_t} = \prod_{i=1}^l \left( \frac{P_i^x}{g^{a'_i \cdot x - e_i} \cdot T_i^{x^2} \cdot S_i} \right)^{k^i} \in \mathbb{G} \quad (10)$$

The verifier can confirm the correctness of all transformations if the prover can prove the knowledge of  $v_t = \sum_{i=1}^l v_i \cdot k^i$  on generator  $h^x / (g^{x^2} u) \in \mathbb{G}$ . To ensure soundness, note that only elements carrying input data ( $P_i$  and  $a'_i$ ) are taking to the first power  $x$ , and other committed values  $S_i, T_i$  are either taking to the second power  $x^2$  or not at all.

Finally, the verifier validates that  $e_i$  doesn't alter the value of  $a_i$ . This can be done by taking the modulus  $q$  of  $e_i$  which must return 0. This is trivial to understand since  $a'_i$  is in  $\mathbb{Z}_q$  so  $e_i$  must be a multiple of  $q$ .

$$\text{if } (e_i \bmod q) \stackrel{?}{=} 0, \text{ then continue} \quad (11)$$



We have so far skipped the overflow problem. If  $a_i + (x \alpha_i \bmod q) > q$ , then we will have an overflow problem in equation 8 10 when computing  $a'_i \cdot x - e_i$ . To get around this, the prover simply needs to check if  $a_i + (x \alpha_i) \bmod q$  overflows  $q$  and subtract  $q \cdot x$  from  $e_i$  if that's the case.

$$\mathbf{if} \ a_i + (x \alpha_i \bmod q) \geq q, \mathbf{then} \ e_i = e_i - x \cdot q \quad i = \{1, \dots, l\} \quad (12)$$

This adjustment does not break the zero-knowledgeness of  $e_i$  since  $x$  (61-bits) is significantly smaller than  $e_i$ 's blinding key  $\omega_i$  (183-bits), so subtracting  $x \cdot q$  from  $e_i$ 's blinding term  $\omega_i \cdot q$  is perfectly unnoticeable except for a negligible chance of  $2^{-122}$ . The validation part of the input-mapping sub-protocol is defined as follows:

*Input* :  $(\vec{P}, \vec{S}, \vec{T} \in \mathbb{G}, \vec{a}' \in \mathbb{Z}_q; \vec{a}, \vec{\alpha}, \vec{\omega} \in \mathbb{Z}_q, \vec{v} \in \mathbb{Z}_p)$   
*P's input* :  $(\vec{P}, \vec{S}, \vec{T}; \vec{a}, \vec{\alpha}, \vec{\omega}, \vec{v})$ ; *V's input* :  $(\vec{P}, \vec{S}, \vec{T})$   
*P compute* :  

$$e_i = ((x \alpha_i \bmod q) - x \alpha_i)x + \omega_i q; \quad i = \{1, \dots, l\}$$

$$\mathbf{if} \ a_i + (x \alpha_i \bmod q) > q,$$

$$\quad \mathbf{then} \ e_i = e_i - q \cdot x \quad i = \{1, \dots, l\}$$
*P*  $\rightarrow$  *V* :  $\vec{e}, \vec{a}'$   
*V*  $\rightarrow$  *P* :  $k \xleftarrow{\$} \mathbb{Z}_p$   
*P compute* :  

$$v_t = \sum_{i=1}^l v_i k^i \in \mathbb{Z}_p$$

$$tr_{v_t} = \text{ProveDL}((h^x / (g^x u)), v_t)$$
*P*  $\rightarrow$  *V* :  $tr_{v_t}$

*V verify inputs* :  

$$\mathbf{if} \ (e_i \bmod q) \stackrel{?}{=} 0, \mathbf{then} \ \text{continue}; \quad i = \{1, \dots, l\}$$

$$\mathbf{else} \ \text{reject}$$

$$PK_{v_t} = \prod_{i=1}^l \left( \frac{P_i^x}{g^{a'_i \cdot x - e_i} \cdot T_i^{x^2} \cdot S_i} \right)^{k^i} \in \mathbb{G}$$

$$\mathbf{if} \ \text{VerifyDL}((h^x / (g^{x^2} u)), PK_{v_t}, tr_{v_t}), \mathbf{then} \ \text{accept}$$

$$\mathbf{else} \ \text{reject}$$

#### Protocol for InputMapping - Verify

**Theorem 1.** *(The Input-Mapping Protocol). The proof system presented in this section has perfect completeness, PHVZK, and CWEE.*

*Proof.* The perfect completeness of protocol InputMapping Validation is trivial to observe.

To prove PHVZK for relation 4, we define a simulator  $\mathcal{S}_{input}$ . To start, simulator  $\mathcal{S}_{input}$  randomly generates  $\vec{S}, \vec{T}$  and sends them to the verifier. After receiving challenge  $x$  from the verifier, the simulator first generates a set of linear polynomials  $\vec{a}'$ , and then simulates the proof transcripts proving the mapping between committed values  $\vec{P}$  and  $\vec{a}'$  it generated.

The simulator  $\mathcal{S}_{input}$  randomly generates and sends  $\vec{e}$  according to the protocol specification ( $e_i$  is generated by first randomly generating a value  $v_i \in \mathbb{Z}_p$  and multiplying it by  $q$  s.t.  $e_i = v_i \cdot q$ ) and

sends them to the verifier. When challenge  $k$  is received, the simulator  $\mathcal{S}_{input}$  calls simulator  $\mathcal{S}_{dlog}$  to generate transcript  $tr_{v_t}$  and send it to the verifier.

The verifier then follows the protocol to compute  $PK_{v_t}$  using transcripts  $\vec{S}, \vec{T}, \vec{a}', \vec{e}$  and calls the *VerifyDL* function to test it against the input transcript  $tr_{v_t}$ . We already know for a fact that there exists a simulator  $\mathcal{S}_{dlog}$  that can simulate witnesses for any discrete-log relation, and that simulators  $\mathcal{S}_{input}$  and  $\mathcal{S}_{dlog}$  choose all proof elements and challenges according to the randomness supplied by the adversary from their respective domains or compute them directly as described in the protocol. Since all elements in proof transcripts are either independently randomly distributed or their relationship is fully defined by the verification equations, we can conclude that protocol InputMapping is PHVZK.

To prove CWEE, we construct an extractor  $\mathcal{X}$  that also uses extractor  $\mathcal{X}_{dlog}$  to extract witnesses from proof of knowledge transcripts (which we know exist for a fact). To start, the extractor  $\mathcal{X}$  interacts with the prover and receives  $\vec{S}, \vec{T}$  from the prover. The extractor  $\mathcal{X}$  then generates a challenge  $x_1$  and forwards it to the prover. After receiving  $\vec{e}_1, \vec{a}'_1$ , the extractor rewinds and repeats this step with another challenge  $x_2$  to retrieve  $\vec{e}_2, \vec{a}'_2$ .

After receiving transcripts  $\vec{S}, \vec{T}$  and transformed inputs  $\vec{a}'$  from the prover, the extractor generates  $k_1$  and then follows the protocol to get  $tr_{v_{t1}}$  (from the prover),  $PK_{v_{t1}}$ . The extractor  $\mathcal{X}$  then calls the extractor  $\mathcal{X}_{dlog}$  to retrieve  $v_{t1}$  from generator  $h^x/(g^{x^2}u)$ . The extractor then rewinds and repeats this step  $l$  times to retrieve  $v_{t2}, \dots, v_{tl+1}$ . Through interpolation, the extractor retrieves witnesses  $v_i$  for all  $i$  in  $\{1, \dots, l\}$ . Since we know for a fact that  $e_i$  cannot alter  $a_i$  and committed values  $\vec{P}, \vec{S}, \vec{T}$  all applied to different powers of  $x$ ,  $a$  cannot be altered except for a negligible probability or we find a non-trivial relationship between generators  $g, h, u$ .

Using any two different challenges  $x_i, x_{i+1}$  we mentioned earlier, the extractor gets  $\vec{a}'_1$  and  $\vec{a}'_2$  from the prover, which we can trivially retrieve  $\vec{a}$  for all  $i = \{1, \dots, l\}$  using the equation below.

$$a'_{1i} - a'_{2i} = \alpha_i(x_1 - x_2) \quad (13)$$

With  $\vec{a}, \vec{\alpha}$  extracted, we can also extract  $\omega$  from equation 7. Plugging witnesses  $\vec{a}, \vec{\alpha}, \vec{v}, \vec{\omega}$  into generators  $g, h, u$ , we can re-write the left and right sides of equation 10 to:

$$(h^x/(g^{x^2}u))^{\sum_i v_i \cdot k_i} = \prod_{i=1}^l \left( \frac{g^{a_i \cdot x} h^{v_i \cdot x}}{g^{a'_i \cdot x - e_i + (v_i - \alpha_i) \cdot x^2 + \omega_i \cdot q} \cdot u^{v_i}} \right)^{k_i} \in \mathbb{G} \quad (14)$$

Take out challenge  $k_i$ , for each  $i$ th input we have:

$$\frac{h^{x \cdot v_i}}{g^{x^2 v_i} u^{v_i}} = \frac{g^{a_i \cdot x} h^{v_i \cdot x}}{g^{a'_i \cdot x - e_i + (v_i - \alpha_i) \cdot x^2 + \omega_i \cdot q} \cdot u^{v_i}} \in \mathbb{G}$$

This implies generator  $g$ 's exponent in the nominator of the right-hand side element must cancel out, which automatically implies the  $g$ 's exponent in the denominator of the right hand side element must be  $a_i \cdot x + x^2 v_i$  for  $g$ 's exponent on the left hand side to be  $x^2 v_i$ . Since we know  $e_i$  cannot alter  $a_i \in \mathbb{Z}_q$  because  $e_i \bmod q = 0$ , we can trivially observe that no other value besides  $a_i$  in  $g$ 's exponent in the denominator ( $(a'_i) \cdot x - e_i + (v_i - \alpha_i) \cdot x^2 + (\omega_i \cdot q)$ ) is multiplied by the single power of  $x$ . This implies the equality above must be true for a computationally bounded prover except for a negligible probability (adversary guessed  $x$  correctly), or we find a non-trivial relationship between generators  $g, h, u$ , and this satisfies our CWEE definition.

## 5 Protocol for Arbitrary Circuits

We first define the relation for the base version of our protocol. For  $l$  input parameters, let  $\mathcal{C}_{\mathbb{F}}$  represent the set of arbitrary arithmetic circuits in  $\mathbb{F}$ , there exists a zero knowledge argument for the relation:

$$\{(g, h, \vec{P}, R \in \mathbb{G}, E \in \mathcal{C}_{\mathbb{F}}; \vec{a}, \vec{v}, r, \phi \in \mathbb{Z}_p) : E(\vec{a}) = r \wedge P_i = g^{a_i} h^{v_i} \wedge R = g^r h^\phi\} \quad (15)$$

$g, h$  are initial public parameters  $pp$  generated during setup. The above relation states that each input parameter to a circuit is represented by a commitment  $P_i$  in  $\mathbb{G}$ , which hides each input value  $a_i$  with a blinding key  $v_i$ .  $r$  is the output of circuit  $E$  computed from inputs  $\vec{a}$ , which is also committed in  $R \in \mathbb{G}$  with blinding key  $\phi$ .

The main idea behind the “input transformation” concept is the process of transforming committed inputs in  $\mathbb{G}$  to linear polynomials in  $\mathbb{F}$  so that the verifier can perform addition and multiplication operations “as is”. For an input commitment  $P_i$  s.t.  $P_i = g^{a_i} h^{v_i} \in \mathbb{G}$  where  $a_i$  is the input value and  $v_i$  is its blinding key, its transformed value linear polynomial is  $a_i' \in \mathbb{Z}_q$  is :

$$a_i' = a_i + x \cdot \alpha_i \in \mathbb{Z}_q \quad (16)$$

$x$  is the challenge provided by the verifier during runtime, and the blinding key of each input is replaced by a random  $\alpha_i$  s.t.  $\alpha_i \neq v_i$ . Likewise, the circuit output commitment  $R = g^r h^\phi \in \mathbb{G}$  also has a matching linear polynomial in  $\mathbb{Z}_q$  with blinding key  $\epsilon$ .

$$r' = r + x \cdot \epsilon \in \mathbb{Z}_q \quad (17)$$

Where  $r'$  is “directly computed” from a list of inputs  $a_i'$  ( $i = \{1, \dots, l\}$ ) by the verifier. Since inputs are in  $\mathbb{Z}_q$ , verifiers can perform arithmetic operations on them just as they do on decrypted numbers. The output value of a circuit evaluation is now a polynomial with  $p_d + 1$  degrees evaluated at point  $x$ . The constant term of the output polynomial is the circuit output  $r$ , and the coefficient of the degree one term is the blinding key  $\epsilon$  of the circuit output.

In the next two sub-sections, we explain our protocol in two steps: In the first step, we introduce a sub-protocol (Protocol InputMapping) that allows the prover to prove each input in  $\mathbb{G}$  is correctly transformed to that in  $\mathbb{Z}_q$ ; In the second step, we introduce the full protocol (Protocol AriCircuit) that uses the aforementioned sub-protocol to validate transformed inputs and proves the circuit output is correctly computed from circuit inputs as relation 15 states. Before jumping into details, there are three key techniques that our protocol leverages that are important to go over:

**Use of two primes** We use two prime domains  $p, q$  in our protocol. All transformed inputs are in field  $\mathbb{Z}_q$ , where the verifier can perform circuit-related computations, and commitments are in  $\mathbb{G}$  of order  $p$ . There are two simple reasons for using  $\mathbb{Z}_q$  in circuit computation: 1) Prime  $p$  is not NTT friendly. 2) It is cheaper to work in a smaller field (we default  $q$  to a 61-bit prime).

**Hiding data in  $\mathbb{Z}_q$  with a significantly larger value in  $\mathbb{F}$**  We use significantly larger numbers in  $\mathbb{F}$  to hide data in  $\mathbb{Z}_q$ . For example, if we want to hide some 61-bit value  $r \in \mathbb{Z}_q$ . The prover can use a 141-bit random value  $\omega$  to hide it s.t.  $r + \omega \in \mathbb{F}$  will perfectly hide  $r$  except for a negligible chance of at least  $2^{-80}$  (e.g., when  $\omega$  is too large or too small, the prover just picks another  $\omega$ ).

The reason we use this hack is because we often perform operations on committed values in  $\mathbb{G}$  with order  $p$ . If the committed value overflows  $p$  after some operation (e.g., for  $r, x \in \mathbb{Z}_p$ , the exponent of  $g^{rx}$  will almost certain overflow  $p$ ), then we cannot use mod  $q$  on the opening of the commitment to get  $rx \in \mathbb{Z}_q$  without changing  $r$  to something else.

**Polynomial degree of a circuit** As mentioned earlier, the degree size of the polynomial generated by our protocol is not necessarily tied to the number of multiplications. In most real-world use cases, we believe the  $p_d$  value generated by our protocol is smaller than the total number of multiplications.

This is because every time an “addition or subtraction” operation is performed, the  $p_d$  value of two input wires merges. For example, if the first input wire  $a_1^7$  has  $p_d = 7$  and the second input wire  $a_1^6$  has  $p_d = 6$ , their sum is an output wire with merged  $p_d = 7$ . The more addition/subtraction operations we have in a circuit, the smaller the  $p_d$  value compared to the total gate count. For example, the sum-of-products circuit explained earlier only has a multiplication depth of 1.

Another technique we use to keep the  $p_d$  value small is a bootstrapping technique that uses “breakers” to reset the  $p_d$  value back to one, a concept we will detail in Section 5.1.

## 5.1 The Protocol

To prove the circuit output is correctly computed from executing transformed inputs  $\vec{a}'$ , the prover provides the proof for all coefficients of the output polynomial. For example, for a simple circuit that computes the sum of two inputs, the prover proves the constant term  $r$  and the coefficient of the degree 1 term  $\epsilon$  of the output polynomial matches the committed values:

$$o = a_1' + a_2' = r + x \cdot \epsilon \quad (18)$$

Computing the output polynomial of the addition operation is the same as adding two polynomials, where  $r = (a_1 + a_2)$  and its blinding key  $\epsilon = (\alpha_1 + \alpha_2)$ . Likewise, multiplying two inputs  $a_1', a_2'$  is the same as multiplying two polynomials:

$$o = a_1' \cdot a_2' = r + x \cdot \epsilon + x^2 \cdot \tau \quad (19)$$

Where  $r = a_1 \cdot a_2$ ,  $\epsilon = a_2\alpha_1 + a_1\alpha_2$ , and  $\tau = \alpha_1 \cdot \alpha_2$ . We use the label “ $o$ ” to represent the circuit output, which is equivalent to the output polynomial evaluated at point  $x$ . The degree of the polynomial ( $p_d$ ) will increase after each multiplication operation, so the efficiency will also drop as there are more computations involved at higher degrees.

To get the linear polynomial  $r'$  from the raw output  $o$ , the verifier needs to subtract out terms with degrees higher than one. In the multiplication circuit above, the verifier needs to eliminate the degree 2 term. To do so, the prover commits to  $\tau$  before the challenge  $x$  is known. When the challenge  $x$  is available, the prover sends  $y$ , the sum of products of coefficients of degree 2 and higher at evaluation point  $x$ , to the verifier and proofs to prove  $x^2\tau = y$ . The verifier can then get the output in the linear polynomial form defined in equation 17:

$$r' = o - y \quad (20)$$

We call each value  $y$  a “breaker”. Breaker(s) subtracts all “noises” (polynomial terms of degree higher than one) from the raw circuit output  $o$ .

In most arithmetic circuits with both addition/subtraction and multiplication/division operations, the circuit output’s polynomial degree may be a lot smaller than the total number of multiplication operations. This is because every time we add two output wires, their polynomials get merged (e.g., circuit  $o = a_1^7 + a_2^3 \cdot a_3^5 + a_4^6$  only outputs a polynomial of degree 8, but 21 multiplications are performed). However, there are special cases where the increase of the  $p_d$  value exceeds the increase in the total number of multiplications.

What we need is 1) a mechanism to allow the prover to repeatedly reset the polynomial degree back to 1 so that the penalty of high degree polynomials can be avoided, and 2) a mechanism to commit to only a sub-linear number of the total  $p_d$  coefficients generated by a circuit.

**Breaking a large circuit into “ $s$ ” smaller circuits** We break the circuit we are evaluating into “ $s$ ” sub-circuits and then batch evaluate them in one run. So when the polynomial degree of a sub-circuit  $i$  reaches degree  $b + 1$ , we reset it to a linear polynomial by evaluating the sub-circuit. The

output of each sub-circuit is then used as the input to the next sub-circuit (Figure 1).

$$\begin{matrix} o_1 \\ o_2 \\ \cdot \\ \cdot \\ o_s \end{matrix} = \begin{pmatrix} r_1 & \epsilon_1 & \tau_{1,1} & \tau_{1,2} & \cdot & \cdot & \tau_{1,b} \\ r_2 & \epsilon_2 & \tau_{2,1} & \tau_{2,2} & \cdot & \cdot & \tau_{2,b} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ r_s & \epsilon_s x & \tau_{s,1} & \tau_{s,2} & \cdot & \cdot & \tau_{s,b} \end{pmatrix} \begin{pmatrix} 1 \\ x \\ \cdot \\ \cdot \\ x^{b+1} \end{pmatrix}$$

Figure 1

We evaluate each sub-circuit  $i = \{1, \dots, s\}$  in the same way we evaluate the full circuit. The raw output of each sub-circuit  $o_i$  is reduced to its linear polynomial form  $r'_i = r_i + x \cdot \epsilon_i$  using breaker  $y_i = \sum_{j=1}^b \tau_{i,j} \cdot x^{j+1}$ .

$$r'_i = o_i - y_i \quad \text{for } i = \{1, \dots, s\} \quad (21)$$

Each  $o_i$  is the output of each sub-circuit at evaluation point  $x$ , and each breaker  $y_i$  is the evaluation output of that sub-polynomial minus the constant term ( $r$ , sub-circuit output) and the degree 1 term ( $\epsilon$ , the blinding key of they sub-circuit output), see Figure 2. Since the output of each sub-circuit is in the same linear polynomial format as inputs to any circuit, they can be reused as inputs to subsequent sub-circuits.

$$\begin{matrix} o_1 \\ o_2 \\ \cdot \\ \cdot \\ o_s \end{matrix} = \underbrace{\begin{pmatrix} r_1 & \epsilon_1 \\ r_2 & \epsilon_2 \\ \cdot & \cdot \\ \cdot & \cdot \\ r_s & \epsilon_s \end{pmatrix}}_{\text{outputs}} \begin{pmatrix} 1 \\ x \end{pmatrix} + \underbrace{\begin{pmatrix} \tau_{1,1} & \tau_{1,2} & \cdot & \cdot & \tau_{1,b} \\ \tau_{2,1} & \tau_{2,2} & \cdot & \cdot & \tau_{2,b} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \tau_{s,1} & \tau_{s,2} & \cdot & \cdot & \tau_{s,b} \end{pmatrix}}_{\text{breakers}} \begin{pmatrix} x^2 \\ x^3 \\ \cdot \\ \cdot \\ x^{b+1} \end{pmatrix}$$

Figure 2

By proving prior knowledge of all coefficients of each sub-circuit (e.g., each row of Figure 2), we know the output  $r'_i$  of each sub-circuit is correct. Sub-circuits are arranged according to the computation order of the full circuit. If outputs of sub-circuits are correct, then the output of the final sub-circuit must also be the output of the full circuit.

It is worth noticing that the prover can set breakers anywhere on a circuit depending on the application design. For simplicity, we assume all breaker are set at  $b + 1$  degrees, where  $b = pa/s$ .

**Make group exponentiation operations sub-linear to the polynomial degree of the full circuit** Using polynomial commitment schemes to evaluate each sub-circuit would be expensive. In our case, the verifiers already know the evaluation output  $o_i$  of each sub-polynomial (sub-circuit) by directly computing them with transformed inputs. Knowing this we can efficiently batch evaluate polynomials by 1) commit and evaluate the output pair  $(r_i, \epsilon_i)$  of each sub-circuit/sub-polynomial, 2) use an intermediary challenge “w” to batch commit reset of the coefficients for all sub-circuits/sub-polynomials.

First, we commit and evaluate the outputs of each sub-circuit using a PCS scheme. Since the prover and verifier runtime cost for running PCS in our protocol is insignificant compared to the circuit size, we pick a group based PCS to minimizing the communication cost. The one we use in our protocol is based on the one defined by Bootle et al. [10]. We define two functions that our protocol will use to commit and evaluate coefficients  $\vec{z}$  at evaluation point  $x$ .

- **polyCommit** $(\vec{c}) \rightarrow \vec{C}$ : commit coefficients  $\vec{c}$  using generators  $\vec{g}$ , returns commitments  $\vec{C} \in \mathbb{G}$  s.t.  $|\vec{C}| = \sqrt{|\vec{c}|}$ .
- **PolyEval** $(\vec{C}, y, x; \vec{c}) \rightarrow b \in \{0, 1\}$ :  $\vec{c}$  are the coefficients of univariate polynomial  $f(X)$  committed in  $\vec{C}$ . At evaluation point  $x$ , this function verifies if  $f(x) = y$ .

Instead of using Pedersen commitments to commit each output as we do with the circuit output  $R$ , we use two polynomial commitments to batch commit the outputs sub-circuits  $\vec{r}, \vec{\epsilon}$ , where  $\vec{r}$  are outputs of circuits and  $\vec{\epsilon}$  are their blinding keys.

$$\beta_r, \beta_\epsilon \xleftarrow{\$} \mathbb{Z}_q \quad (22)$$

$$\text{polyCommit}(\beta_r || \vec{r}) \rightarrow: \vec{R}_s \quad (23)$$

$$\text{polyCommit}(\beta_\epsilon || \vec{\epsilon}) \rightarrow: \vec{E}_s \quad (24)$$

To avoid possible leaks of data, we add a blinding output  $r_\beta$  and corresponding blinding key  $r_\epsilon$  to the commitments.

When challenge  $x$  is known, the prover sends  $r'_1, \dots, r'_s$  to the verifier. The verifier uses the PCS scheme at some evaluation point  $z$  ( $y_r, y_\epsilon$ ) and the following equality to check if the transcripts  $\vec{r}'$  match those committed earlier.

$$\text{polyEval}(\vec{R}_s, y_r, z; \vec{r}) \wedge \text{polyEval}(\vec{E}_s, y_\epsilon, z; \vec{\epsilon}) \quad (25)$$

$$0 = (y_r + y_\epsilon \cdot x - (r'_\beta + \sum_{i=1}^s r'_i)) \bmod q \quad (26)$$

The output  $o_i$  of each polynomial evaluated at point  $x$  is directly computed by the verifier using transformed inputs. Since  $r'_i$  is a committed value, we can infer that each breaker  $y_i$  is equivalent to a committed value because it can be directly computed from  $o_i$  and  $r'_i$ .

$$y_i = o_i - r'_i \quad \text{for} \quad i = \{1, \dots, s\} \quad (27)$$

Second, we use an intermediary challenge “ $w$ ” to “batch” commit each  $j$ -th term as a single coefficient without breaking the soundness of the protocol. This is possible because the constant term  $\vec{r}$  is already committed, and we know  $\vec{\sigma}$  are correct.

The challenge  $w$  is made available after the prover commits to  $\vec{r}, \vec{\epsilon}$  but before the challenge  $x$  is available. The “batched”  $j$ -th coefficient is defined as:

$$c_j = \sum_{i=1}^s \tau_{i,j} \cdot w^i \in \mathbb{Z}_q \quad j = \{0, \dots, b\} \quad (28)$$

$c_0$  ( $j = 0$ ) is a set of blinding keys  $\tau_{1,0}, \dots, \tau_{s,0}$  we use to hide blinding keys  $\epsilon_i$  to cover cases when  $\epsilon_i = 0$  (which may happen in a boolean circuit).

$$\epsilon_i = \epsilon_i - \tau_{i,0} \in \mathbb{Z}_q \quad (29)$$

The updated matrix with blinding elements  $\tau_{1,0}, \dots, \tau_{s,0}$  is presented in Figure 3 below.

$$\begin{matrix} y_1 w \\ y_2 w^2 \\ \cdot \\ \cdot \\ y_s w^s \end{matrix} = \begin{pmatrix} \tau_{1,0} w & \tau_{1,1} w & \tau_{1,2} w & \cdot & \cdot & \tau_{1,b} w \\ \tau_{2,0} w^2 & \tau_{2,1} w^2 & \tau_{2,2} w^2 & \cdot & \cdot & \tau_{2,b} w^2 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \tau_{s,0} w^s & \tau_{s,1} w^s & \tau_{s,2} w^s & \cdot & \cdot & \tau_{s,b} w^s \end{pmatrix} \begin{pmatrix} x \\ x^2 \\ \cdot \\ \cdot \\ x^{b+1} \end{pmatrix}$$

Figure 3

When challenge  $x$  is available, the verifier validates  $\vec{y}$  by multiplying each  $y_i$  by  $w^i$  and each  $c_j$  by  $x^{j+1}$ . The equality below must be true.

$$\sum_{j=0}^b c_j \cdot x^{j+1} = \sum_{i=0}^s y_i \cdot w^i \in \mathbb{Z}_q \quad (30)$$

To show why this is sound for all sub-circuit outputs. Let's say  $r_i^* = r_i - \delta_i$  for some arbitrary  $\delta_i$ , the dishonest prover needs to find a set  $\vec{\Delta}$  before challenge  $x$  is known s.t.

$$\sum_{j=0}^b (c_j + \Delta_j) \cdot x^{j+1} = \sum_{i=0}^s (y_i + \delta_i) \cdot w^i \in \mathbb{Z}_q$$

The equality above can be rewritten as the equality below, which clearly shows such  $\vec{\Delta}$  cannot be found without prior knowledge of challenge  $x$ .

$$\sum_{j=0}^b \Delta_j \cdot x^{j+1} = \sum_{i=0}^s \delta_i \cdot w^i \in \mathbb{Z}_q$$

Lastly, sending  $\vec{c}$  to the verifier may not be safe, the more secure and communication-cost-efficient option is to use a polynomial commitment scheme to commit and evaluate  $\vec{c}$  at evaluation point  $x$ .

$$\text{polyCommit}(\vec{c}) \rightarrow: \vec{C} \quad (31)$$

Let  $y_c \in \mathbb{Z}_q$  be the evaluation output at point  $x$  for a univariate polynomial with coefficients  $\vec{c}$  committed in  $\vec{C}$ . The equation 30 is now updated to:

$$\text{polyEval}(\vec{C}, y_c, x; \vec{c}) \wedge y_c = \left( \sum_{i=1}^s y_i \cdot w^i \right) \bmod q \quad (32)$$

We now have all the necessary pieces to describe our main protocol for arithmetic circuits.

## 5.2 The Complete Protocol For Arithmetic Circuits

We define two more functions for our protocol definition. function

**computeSubCircuitKeys** is used by the prover to compute the keys of each sub-circuit or “row” in Figure 3, and function

**computeSubCircuit** is used by the verifier to compute the value of a sub-circuit at the evaluation point  $x$ :

1. function **computeSubCircuitKeys<sub>i</sub>**(“input values”, “input keys”): for  $i = \{1, \dots, s\}$ , it takes input values  $\vec{a}$  and keys  $\vec{\alpha}$  to evaluate the sub-circuit and outputs  $r_i, \epsilon_i, \vec{\tau}_i$  (coefficients of  $o_i$ ).
2. function **computeSubCircuit<sub>i</sub>**(“inputs in linear polynomial form”, “output from the previous computeSubCircuit function”): for  $i = \{1, \dots, s\}$ , it trivially computes the result  $o_i$  from the inputs to the sub-circuit.

For example, if the logic of the  $i$ th sub-circuit is to return the product of  $l$  inputs, then the **computeSubCircuit<sub>i</sub>** function simply performs  $o_i = a_1 \times a_2 \times \dots \times a_l$ . Since  $a_1, \dots, a_l$  are linear polynomials evaluated at point  $x$ ,  $o_i$  is the evaluation of the output polynomial at point  $x$ , and the **computeSubCircuitKeys<sub>i</sub>** function computes all coefficients of the output polynomial. We are now ready to introduce the complete version of our protocol - Protocol AriCircuit - as follows:

$$\text{Input} : (\vec{P}, R \in \mathbb{G}; \vec{a}, \vec{v}, r, \phi \in \mathbb{Z}_p) \quad (33)$$

$$\mathcal{P}'s \text{ input} : (\vec{P}, R; \vec{a}, \vec{v}, r, \phi); \quad \mathcal{V}'s \text{ input} : (\vec{P}, R) \quad (34)$$

$$\mathcal{P} \text{ compute} : \quad (35)$$

$$\alpha_i \xleftarrow{\$} \mathbb{Z}_q, \quad i = \{1, \dots, l\} \quad (36)$$

$$\beta_r, \beta_\epsilon \xleftarrow{\$} \mathbb{Z}_q \quad (37)$$

$$\mathbf{for} \quad i = 1, \dots, s \quad \{ \tag{38}$$

$$r_i, \epsilon_i, \vec{\tau}_i = \text{computeSubCircuitKeys}_i(\vec{a}', \vec{a}', r_i, \epsilon_i, \vec{\tau}_i); \tag{39}$$

$$\epsilon_i = (\epsilon_i - \tau_{i,0}) \bmod q \quad \} \tag{40}$$

$$\text{polyCommit}(\beta_r || \vec{r}) \rightarrow: \vec{R}_s; \quad \text{polyCommit}(\beta_\epsilon || \vec{\epsilon}) \rightarrow: \vec{E}_s \tag{41}$$

$$\text{InputMapping-Setup}(\vec{a} || r_s, \vec{a} || \epsilon_s, \vec{v} || \phi) \rightarrow: \vec{S}, \vec{T} \tag{42}$$

$$\mathcal{P} \rightarrow \mathcal{V} : \vec{R}_s, \vec{E}_s, \vec{S}, \vec{T} \tag{43}$$

$$\mathcal{V} \rightarrow \mathcal{P} : w \xleftarrow{\$} \mathbb{Z}_p \tag{44}$$

$$\mathcal{P} \text{ compute :} \tag{45}$$

$$\tau_{i,0} \xleftarrow{\$} \mathbb{Z}_q, \quad i = \{1, \dots, s\} \tag{46}$$

$$c_j = \sum_{i=1}^s (\tau_{i,j} \cdot w^i) \in \mathbb{Z}_q \quad j = \{0, \dots, b\} \tag{47}$$

$$\text{polyCommit}(\vec{c}) \rightarrow: \vec{C} \tag{48}$$

$$\mathcal{P} \rightarrow \mathcal{V} : \vec{C} \tag{49}$$

$$\mathcal{V} \rightarrow \mathcal{P} : x \xleftarrow{\$} \mathbb{Z}_p \tag{50}$$

$$\mathcal{P} \text{ compute :} \tag{51}$$

$$a'_i = a_i + x \cdot \alpha_i \in \mathbb{Z}_q \quad i = \{1, \dots, l\} \tag{52}$$

$$r'_i = r_i + x \cdot \epsilon_i \in \mathbb{Z}_q \quad i = \{1, \dots, s\} \tag{53}$$

$$r'_\beta = \beta_r + x \cdot \beta_\epsilon \in \mathbb{Z}_q \tag{54}$$

$$y_c = c_0 + \sum_{j=1}^b c_j \cdot x^j \in \mathbb{Z}_q \tag{55}$$

$$\mathcal{P} \rightarrow \mathcal{V} : \vec{a}', \vec{r}', r'_f, y_c \tag{56}$$

$$\mathcal{V} \rightarrow \mathcal{P} : z \xleftarrow{\$} \mathbb{Z}_q \tag{57}$$

$$\mathcal{P} \text{ compute :} \tag{58}$$

$$y_r = \beta_r + \sum_{i=1}^s r_i \cdot z^i \in \mathbb{Z}_q; \quad y_\epsilon = \beta_\epsilon + \sum_{i=1}^s \epsilon_i \cdot z^i \in \mathbb{Z}_q \tag{59}$$

$$\mathcal{P} \rightarrow \mathcal{V} : y_r, y_\epsilon \tag{60}$$

$$\mathcal{V} \text{ verify final output :} \tag{61}$$

$$\mathbf{for} \quad i = 1, \dots, s \quad \{ \tag{62}$$

$$o_i = \text{computeSubCircuit}_i(\vec{a}' || \vec{r}') \in \mathbb{Z}_p \tag{63}$$

$$y_i = o_i - r'_i \in \mathbb{Z}_q \quad \} \tag{64}$$

$$\mathbf{if} \quad (\text{polyEval}(\vec{R}_s, y_r, z; \vec{r}) \wedge \text{polyEval}(\vec{E}_s, y_\epsilon, z; \vec{\epsilon})) \tag{65}$$

$$\wedge 0 = (y_r + y_\epsilon \cdot x - (r'_\beta + \sum_{i=1}^s r'_i)) \bmod q \tag{66}$$

$$\wedge \text{polyEval}(\vec{C}, y_c, x; \vec{c}) \tag{67}$$

$$\wedge y_c = (\sum_{i=1}^s y_i \cdot w^i) \bmod q \quad \mathbf{then} \text{ continue} \tag{68}$$

$$\mathbf{else} \text{ reject} \tag{69}$$



if InputMapping-Verify( $\vec{P}||R, \vec{S}, \vec{T}, \vec{a}'||r'_s; \vec{a}||r_s, \vec{\alpha}||\epsilon_s, \vec{v}||\phi$ ) (70)

then continue (71)

else reject (72)

### Protocol AriCircuit

**Theorem 2.** (Protocol AriCircuit). *The proof system presented in this section has perfect completeness, PHVZK, and CWEE.*

*Proof.* The perfect completeness of protocol AriCircuit is trivial to observe.

To prove PHVZK for relation 15, we define a simulator  $\mathcal{S}$ . Simulator  $\mathcal{S}$  calls on simulators  $\mathcal{S}_{input}$  defined earlier to simulate interactions in the InputMapping sub-protocol used in AriCircuit, and uses Simulator  $\mathcal{S}_P$ , which we know exist for the PCS scheme we use, to simulate interactions in the polynomial commitment scheme our protocol calls.

We have already shown that  $\mathcal{S}_{input}$  can simulate all interactions needed in sub-protocol InputMapping. We now show how  $\mathcal{S}$  generates the rest of the transcripts according to the randomness supplied by the adversary from their respective domains or computes them directly as described in the protocol.

Simulator  $\mathcal{S}$  randomly generates input witnesses  $\vec{a}^*, \vec{a}'^*$  and computes circuit output witnesses  $\vec{r}^*, \vec{e}^*, \vec{\tau}^*$  and commits them to  $\vec{R}_s, \vec{E}_s, \vec{C}$  according to the protocol specification and sends them to the verifier. The simulator then follows the protocol to create  $\vec{c}^*$  using challenge  $w$  received from the verifier and sends them to the verifier. After receiving challenge  $x$  from the verifier, the simulator computes  $\vec{a}'^*, \vec{r}'^*$  according to the protocol specification with generated input witnesses and sends them to the verifier. The simulator  $\mathcal{S}$  follows the protocol specification to create all other transcripts needed to complete the protocol, and rewindings only take place inside the simulator  $\mathcal{S}_{input}$  (for simulating transcripts in InputMapping-Verify function) and simulator  $\mathcal{S}_P$  (for simulating transcripts in the polyEval function). This implies that if the input-mapping process and PCS scheme are PHVZK, then our protocol is PHVZK.

Since all elements in proof transcripts are either independently randomly distributed or their relationship is fully defined by the verification equations, we can conclude that the protocol AriCircuit is PHVZK.

To prove CWEE, we define extractor  $\mathcal{X}$  that calls on extractors  $\mathcal{X}_{input}$  defined earlier to extract witnesses for the sub-protocol used in AriCircuit. We also define an extractor  $\mathcal{X}_P$  that extracts witnesses from a polynomial commitment, which is defined for the PCS scheme we use.

We already know  $\mathcal{X}$  can extract  $\vec{a}, \vec{\alpha}$  and  $\vec{r}, \vec{e}$  using extractor  $\mathcal{X}_{input}$  from committed transcripts  $\vec{S}, \vec{T}$ . Using extracted input witnesses, we can call the function computeSubCircuitKeys to compute circuit witnesses  $\vec{r}, \vec{e}, \vec{\tau}$ . In the next few paragraphs, we show these circuit witnesses must match those extracted from circuit transcripts  $\vec{r}', \vec{R}_s, \vec{E}_s, \vec{C}$ .

First, the extractor  $\mathcal{X}_P$  extracts coefficients  $\vec{c}$  from polynomial commitments  $\vec{C}$ . With  $\vec{c}$ , the extractor  $\mathcal{X}$  rewinds one more time to line 43 and then rewinds another  $s$  times generates  $s+1$  challenges  $w_1, w_2, \dots, w_{s+1}$ , using which the extractor can extract witnesses  $\vec{\tau}_j$  for  $j \in \{0, \dots, s\}$  using interpolation. Rearrange equality 30 and substitute  $\vec{c}$  for  $\vec{\tau}, w$  as specified in 28 we get the following equality:

$$\sum_{j=1}^b \left( \sum_{i=1}^s \tau_{i,j} w^i \right) \cdot x^{j+1} = \sum_i y_i \cdot w^i \in \mathbb{Z}_q \quad (73)$$

If  $y_1, \dots, y_s$  on the right hand side the equality above is correct, then the equality above can only be true for any challenge  $x$  provided each  $z_j$  is also correctly computed using extracted  $\vec{\tau}_j$  with equation 28 except for a negligible probability. To check if  $y_1, \dots, y_s$  are correct, we check the soundness of  $r'_\beta, r'_1, \dots, r'_s$ .

Second, the extractor  $\mathcal{X}_P$  extracts coefficients  $\vec{r}, \vec{\epsilon}$  from polynomial commitments  $\vec{R}_s, \vec{E}_s$  using evaluation points  $\vec{z}$ . With  $\vec{r}, \vec{\epsilon}$  extracted, the equality of equation 26 can be rewritten to:

$$r'_\beta + \sum_{i=1}^s r'_i \cdot z^i = \beta_r + \beta_\epsilon \cdot x + \sum_{i=1}^s (r_i + \epsilon_i \cdot x) \cdot z^i \in \mathbb{Z}_q \quad (74)$$

The equality above can only be true for any challenge  $x$  at any evaluation point  $z$  if transcripts  $r'_\beta$  and each  $r'_i$  (left hand side of the equality) are computed according to equations 16, 17 using committed witnesses  $\vec{r}, \vec{\epsilon}$  shown on the right hand side of the equality above.

Third, using the extracted  $\vec{r}, \vec{\epsilon}$ , we update the equality 73 to the following:

$$\sum_{j=1}^b \left( \sum_{i=1}^s \tau_{i,j} w^i \right) \cdot x^{j+1} = \sum_i^s (o_i - (r_i + x \cdot \epsilon_i)) \cdot w^i \in \mathbb{Z}_q \quad (75)$$

Knowing that  $\vec{r}, \vec{\epsilon}$  are committed witnesses and that  $\vec{o}$  must be correct since it was directly computed by the verifier, the equality above can only be true for any pairs of challenges  $w, x$  if coefficients for each sub-circuit  $\vec{r}, \vec{\epsilon}, \vec{o}$  are correct except for a negligible probability.

Finally, we check if the extracted circuit witnesses  $\vec{r}, \vec{\epsilon}, \vec{o}$  extracted from circuit transcripts match those computed from input witnesses  $\vec{a}, \vec{\alpha}$  using computeSubCircuitKeys functions. Since the coefficients computed from computeSubCircuitKeys functions also need to satisfy equalities 73, 74, 75 for the same evaluation result  $\vec{o}$  for any pair of challenges  $w, x$ . The witnesses computed from  $\vec{a}, \vec{\alpha}$  must match those extracted from circuit transcripts except for a negligible probability, or else we find a non-trivial discrete log relationship between generators  $g, h$  (for input witnesses).

### 5.3 Customized Sub-Circuit and The Use of Range Proof for Comparison Operations

One of the primary reasons for using a boolean circuit over an arithmetic circuit in the real world (there are no real-world applications of trying to prove a hash) is the ability to perform comparison operations ( $>, <, \geq, \leq, =$ ). Our design allows the use of customized circuit(s) to embed range-proof protocols to evaluate comparison operations inside the arithmetic circuit being processed. This way, there will be fewer needs for expensive boolean circuits and/or the expensive process of decomposing/recomposing integers to bits within a circuit.

The idea is similar to that of “custom gates” found in SNARKs protocols in principle but very different in implementation. In general, the design goal of a customized circuit is to utilize existing algorithms/protocols to handle operations that would otherwise be expensive in our protocol (or any other protocol).

In particular, we can use range proof inside an arithmetic circuit to handle all comparison operations and decimal point reductions. This is huge in practice because either using a boolean circuit directly or converting to/from a boolean circuit inside an arithmetic circuit is expensive.

For example, to prove  $a_1 > a_2$  (or  $P_1 > P_2$ ), the prover can do the following:

1. Commit to their difference  $C = g^c h^v$  s.t.  $c = a_1 - a_2$  (or compute  $C$  from  $P_1, P_2$  using additive homomorphism e.g.  $C = P_1/P_2$ ).
2. Call protocol InputMapping to check  $c' = a'_1 - a'_2 \in \mathbb{Z}_q$  maps  $C \in \mathbb{G}$ ;
3. Use a range-proof protocol to prove  $C > 0$ . If returns true, then we know  $a_1 > a_2$ .

An example usage is as follows:

```

computeSubCircuit : ( $\vec{a}' \in \mathbb{Z}^q, \vec{C} \in \mathbb{G}$ )
   $c'_1 = a'_1 - a'_2 \in \mathbb{Z}_q$ 
  if Protocol RangeProof( $C_1, 0$ )
     $c'_2 = a'_3 - a'_4 \in \mathbb{Z}_q$ 
    if Protocol RangeProof( $C_2, 0$ )
       $o = \text{do something}$ 
    else
       $o = \text{do something else}$ 
  else
     $o = \text{do something}$ 
  if Protocol InputMapping( $C_1 || C_2, c'_1 || c'_2$ ) return  $o$ 
  else reject

```

Function computeSubCircuit (Customized)

In the computeSubCircuit function defined above, the circuit first tests if  $a'_1 > a'_2$  and then tests if  $a'_3 > a'_4$ . Before the function returns  $o$ , it batch checks the mapping between each  $c'_i$  and  $C_i$  pair. In practice, all calls to range proof should also be batch verified at the end of the function.

We can bypass the “inactive” part of the circuit (similar to that of suBlonk [18]). For example, if the first range proof returns false (e.g.,  $a_1 < a_2$ ), then both the prover and the verifier can bypass the two “else” parts of the circuit above. However, it is worth noticing that using a customized sub-circuit bypassing parts of the circuit may leak information about data to attackers, so one must use such a strategy with extreme caution.

We believe combining the arithmetic circuit and range-proof protocols is the most efficient way to run a zero-knowledge test in the real world. This is much more powerful than it seems. Besides handling comparison operators, one can also use embedded range proof to verify floating point operations (perform multiplication/division operations as full integer operations, then remove decimal places by proving their range). We believe it will allow us to use arithmetic circuits to handle almost all types of business logic that would otherwise require boolean circuits.

It is also not necessary that all breakers  $y_1, \dots, y_m$  have the same  $b$  value. For example, if some value  $a_i$  is taken to its 100th power ( $a_i^{100}$ , degree term  $s = 100$ ) and will be used as inputs in multiple places of a circuit, then it would be wise to use a breaker to cut its degree to 1 (in linear polynomial form  $a_i^{100} + X\alpha_i$ ) before being used as inputs in other places of a circuit.

#### 5.4 Memory Efficiency

The memory consumption cost of our protocol is  $O(s)$ . However, our design approach allows us to improve the memory consumption cost of our protocol to  $O(b + 2s)$ . We only need to do two thing to achieve a memory cost of  $O(b + 2s)$ : First, delete  $\vec{\tau}$  from memory in line 42; Second, recompute  $\vec{\tau}_i$  after challenge  $x$  is available so that  $y_t$  can be computed ( in line 55 ).

```

for  $i = 1, \dots, m$  {
   $r_i, \epsilon_i, \vec{\tau}_i = \text{computeSubCircuitKeys}_i(\vec{a}', \vec{\alpha}')$ ;
   $\epsilon_i = (\epsilon_i - \mu_i) \bmod q \in \mathbb{Z}_q$ ;
   $\vec{a}' = \vec{a} || r_i, \quad \vec{\alpha}' = \vec{\alpha} || \epsilon_i$ ;

   $z_j = z_j + \sum_{i=1}^s (\tau_{i,j} \cdot w^i) \in \mathbb{Z}_q \quad j = \{1, \dots, b\}$ ;

```

}

Since we have  $b = s = \sqrt{p_d}$  in the default setting, the asymptotic memory cost is  $O(\sqrt{p_d})$ .

## 5.5 Cost Analysis

The prover runtime of our protocol is dominated by  $O(\iota s \log s + l)$  field operations and  $O(b + l)$  group exponentiations. We set  $b = \sqrt{s}$  in our benchmark testing, so the cost of group exponentiations grow sub-linearly to the total polynomial degree generated by the circuit. The value of  $\iota$  depends on how the circuit is wired. For the sequential multiplication test case that we benchmark against,  $\iota = m \sum_{i=1}^{\log b} i$ ; the verifier runtime is dominated by  $O(n + l)$  field operations and  $O(m + b + l)$  group exponentiations; and the communication cost is dominated by  $O(b + l)$  group elements and  $O(m + l)$  field elements.

Our protocol is also natively faster than its asymptotic cost indicates because group exp. operations of our protocol operate mostly in  $q$  (61-bits), which is significantly smaller than  $p$  in ECC. This gives us approximately 2.5X performance gains when performing multi-exponentiations over standard ECC multiplications. Although the verifier runtime is technically linear, it is so efficient to the point that it is close to SNARKs with trusted setup. This is because all the asymptotically slow operations are performed at field level (many papers consider this free).

It is important to note that  $p_d \neq n$ . For some arithmetic circuits, the total polynomial degree ( $p_d$ ) can be smaller than the number of multiplications because every time we perform the “addition” operation on two output wires. On the other hand, if  $p_d$  is bigger than the total number of multiplications, then we may need to set breakers more aggressively to cut the  $p_d$  value down to an acceptable level, preferably lower than  $n$ .

It is also worth noting that it is not necessary to set all breakers at some fixed point  $b$ . For example, if we know several output wires with high degrees are going to get merged through addition operations, we may want to save the breaker until they are merged together.

## 6 Performance Comparison

We compare the performance of our protocol to some of the most popular transparent zero-knowledge protocols for which open source codes are available. Our test runs are performed on an Intel(R) Core(TM) i7-9750H CPU @ 2.60 Ghz. All tests are run on a single CPU thread. Our test code is a non-interactive implementation (using Fiat-Shamir heuristic). For group operations, we use the curve25519-dalek implementation, and Pippenger acceleration is applied to all sum-of-product group operations. For field operations, we use the Montgomery algorithm to accelerate modular multiplications on the prime  $q$ .

We compare our protocol against Hyrax, Liger, Aurora, and Spartan-NIZK. These protocols were chosen because they are the most representative of popular zero-knowledge protocols and can be verified with open source code. In particular, Aurora outperforms STARK in all key parameters (prover/verifier runtime, proof size), and the NIZK version of Spartan offers the most balanced performance across all performance parameters. We do not consider SNARKs because they are hardly efficient after switching to transparent mode.

We didn’t consider protocols that depend on circuit depth, such as GKR-based protocols, because they cannot handle  $2^{20}$  sequential multiplications.

Spartan++ and Lakonia are two more recent developments that we didn’t include in our benchmark testing but are worth mentioning. The improvement of Spartan++ over SpartanNIZK is marginal, and the performance of Lakonia is largely comparable to that of SpartanNIZK. (The prover performance of SpartanNIZK is approximately 3X more efficient, and the verifier performance is 1.5X more efficient than that of Lakonia, while Lakonia is 4X more efficient than SpartanNIZK in proof size).

We set the number of inputs to our protocol to 20 integers. The circuit we use performs  $n$  sequential multiplications on  $l$  inputs. This is because we want to demonstrate that our protocol can handle

high-depth circuits (e.g.,  $p_d = n$ ). If we run a shallow circuit where polynomial degree is small, the benchmark result will likely be significantly better. For example, if we have a circuit that computes the sum-of-products of inputs where  $n = 2^{20}$  (e.g.,  $\sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 \dots + a_n b_n$ ), the prover will complete proving such a shallow circuit (one million multiplications and one million sums) in under 50 milliseconds for such a shallow circuit with  $p_d = 1$  (group exponentiation cost is approaching zero and NTT is not even necessary for field operations).

There are special cases where in some sub-circuits we have  $p_d$  value grow faster than the number of multiplications performed (e.g.,  $o = (((a'_1)^2)^2)^2$ ). In such cases, we may want to place breakers more aggressively to better handle the situation. This can be done by either setting fixed  $b < \sqrt{p_d}$  or placing more breakers inside these uncommon sub-circuits. In a real-world situation where there are both multiplication and addition operations, it is unlikely that  $p_d$  is bigger than  $n$ .

To maximize the advantage of the NTT algorithm in computing sequential multiplications, we process each segment ( $subCircuit_i$ ) of our circuit in binary tree format to represent layers we would see in the real world. Such tuning will likely not be required in real-world applications since large circuits are usually layered and multiplication gates should be somewhat balanced out across layers already.

## 6.1 Benchmark Results

We set  $b = s = \sqrt{p_d}$  to get a more balanced result. Alternatively, one can set  $b$  smaller to get better prover runtime performance in exchange for more expensive verifier and communication costs (section 5.2). This is because doing so will 1) may significantly cut down the polynomial degree  $p_d$  value of the circuit. 2) evade expensive NTT computations at high degrees and better leverage Pippenger acceleration in computing  $R_s, E_t$ , which will continuously improve group exponentiation operations before peaking out at around  $m = 2^{14}$ .

**Table 1.** Prover performance comparison (seconds)

Circuit size	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$
Hyrax	1	2.8	9	36	117	486
Ligero	0.1	0.4	1.6	4	17	69
Aurora	0.5	1.6	6.5	27	116	485
SpartanNIZK	0.02	0.05	0.16	0.6	1.7	6.2
This work	0.003	0.005	0.01	0.04	0.16	0.7

Table 1 shows that as the circuit size gets bigger, the prover performance of our protocol is becoming increasingly more efficient than all the other protocols we are comparing against. This is because the cost associated with the number of inputs to the circuit is fixed (20 inputs), and its impact relative to the cost of evaluating the whole circuit gradually declines as the circuit size gets bigger (the same effect will also apply to verifier runtime and proof size benchmarks below). To the best of our knowledge, our protocol offers the best prover performance in the non-interactive setting in the literature.

Table 2 shows that the communication cost of our protocol dominates that of Ligero and Aurora, while being largely comparable to SpartanNIZK and Hyrax. The fixed cost of one additional input is 112 bytes, you can add or subtract as many inputs as needed to get the communication cost that fits your scenario rather than take the default  $l = 20$ . Please note that other protocols also incur comparable costs when they map more witnesses to some pre-committed/encrypted value in the public data store.

Table 3 demonstrates that our protocol achieves a significant improvement of over one order of magnitude in verifier runtime compared to other protocols.

**Table 2.** Proof size comparison (kilobytes)

Circuit size	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$
Hyrax	14	17	21	28	38	58
Ligero	546	1,076	2,100	5,788	10,527	19,828
Aurora	477	610	810	1,069	1,315	1,603
SpartanNIZK	9	12	15	21	30	48
This work	3.3	4	5.3	7	10.5	16.2

**Table 3.** Verifier performance comparison (milliseconds)

Circuit size	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$
Hyrax	206	253	331	594	1.6s	8.1s
Ligero	50	179	700	2s	7.5s	33s
Aurora	192	590	2s	7.2s	29.8s	118s
SpartanNIZK	7	11	17	36	103	387
This work	1.4	1.5	1.7	2.0	3.3	6.8

The NTT friendly prime number  $q$  we used for our benchmark testing is 1945555039024054273, a 61-bit prime that implies the soundness error will be at most  $\frac{b}{q} = \frac{b}{s} = 2^{-51}$  in our test case ( $b = s = \sqrt{pd}$ ), which is good enough for many applications and ones where one interaction is allowed.

When a soundness error of  $2^{-51}$  is not enough, the dumb and straight-forward way is to run the whole protocol twice to get a soundness error of at least  $(\frac{b}{q})^2 = 2^{-102}$ . This will almost double the cost of everything (the prover only needs to compute vector commitment  $R_s$  once), but our protocol will still claim the title of the fastest prover and verifier runtime in the literature by a wide margin. The more advanced way is to use a bigger  $q$  prime. For example, a 90-bit  $q$  prime will comfortably increase the soundness error to at least  $2^{-80}$ , just that there would be a lot of engineering work to get an efficient NTT and modular arithmetic implementation at a higher bit value. e.g., at 128-bit, which will require optimization at the assembly level for which no open source code is available at the moment, unlike that for 64-bit (come with the CPU) and 256-bit (optimized over the years because of ECC implementations).

It is worth noticing that input transformation costs can be shared across multiple circuits if the inputs are reused as inputs to other circuits, which may lead to further reductions in communication costs in the real world.

## 6.2 Benchmark at different $b$ value

While our protocol can be tailored to support boolean circuits, it is optimized for arithmetic circuits that can handle programming logics directly. To get an efficient prover for a boolean circuit, we probably need to set breaker more aggressively. Table 4 shows the performance number when running a circuit with  $n = 2^{20}$  multiplication operations at different ‘ $b$ ’ value in the standard mode.

The prover runtime performance peaks at around 0.5 seconds when  $b = 2^6$ , then it starts to increase again afterwards. This is because the cost of computing the polynomial commitments  $R_s, E_s$  gets increasingly expensive as the number of sub-circuits ( $m$ ) increases ( $< 10\%$  when  $b = 2^{10}$ , and  $> 90\%$  when  $b = 2^3$ ). Therefore, the cost of performing group exponentiation operations committing  $\vec{r}, \vec{e}$  is getting more expensive relative to the shrinking cost of field operations computing  $\vec{r}, \vec{e}, \vec{\tau}$  as  $b$  decreases. One thing worth notice is that the soundness error  $\frac{b}{s}$  will increase as  $s$  gets increasingly bigger than  $\sqrt{pd}$ .

**Table 4.** Performance comparison at different ‘b’ value ( $n = 2^{20}$ )

b	Prover	Verifier	Communication Cost
$2^{10}$	0.70 s	6.8 ms	16 KB
$2^9$	0.63 s	7.2 ms	26 KB
$2^8$	0.56 s	8.5 ms	43 KB
$2^7$	0.51 s	10 ms	79 KB
$2^6$	0.50 s	12 ms	147 KB
$2^5$	0.56 s	14 ms	283 KB
$2^4$	0.72 s	19 ms	547 KB
$2^3$	1.08 s	38 ms	1075 KB

The down side of lowering ‘b’ value is that the communication cost will gets higher for lower ‘b’ value. Fortunately, boolean circuits usually have hundreds of bits running in parallel, and their outputs are also in hundreds of bits, so the true cost of circuit polynomial degree will also get evened out across these bits.

The other good news is that table 4 is a little bit deceiving because the smaller the  $b$  may also lead to a smaller polynomial degree of a circuit (e.g., when some output of a sub-circuit with a high degree is used multiple times to multiply other values of a circuit.), which in turn leads to faster runtime performance.

Regardless, we believe the most intuitive way to use zero knowledge proof for business users and developers is to avoid using boolean circuit altogether, unless it is for specialized use cases like proof knowledge of a hash’s pre-image.

### 6.3 Benchmark at different $b$ value in the memory efficient setting

One of the biggest advantages of our protocol is that it provides a blueprint for developing a fast, memory-efficient, non-interactive zero-knowledge protocol. The theoretical memory cost is  $O(b + s)$ . Using the memory-efficient implementation mentioned in Section 4.5, we get the results listed in Table 4.

**Table 5.** Performance comparison in memory efficient setting ( $n = 2^{20}$ )

b	Prover	Verifier	Communication Cost
$2^{10}$	1.36 s	6.8 ms	16 KB
$2^9$	1.3 s	7.2 ms	26 KB
$2^8$	1.1 s	8.5 ms	43 KB
$2^7$	0.95 s	10 ms	79 KB
$2^6$	0.92 s	12 ms	147 KB
$2^5$	0.95 s	14 ms	283 KB
$2^4$	1.11 s	19 ms	547 KB
$2^3$	1.59 s	38 ms	1075 KB

The memory efficient mode peaks at around 1.08 million multiplications per second when  $b = 2^6$ . Since the number of NTT operations is doubles, the performance is noticeably worse than that shown in Table 4 running in the non-memory efficient mode.

Regardless, the benchmark numbers shown in Table 4 compare well against top-of-the line VOLE-based protocols (shown in Table 5; reported numbers are copied directly from their paper [38]), given

that our protocol is non-interactive and offers a significantly smaller proof size without requiring pre-setup like that of Ant-Man.

**Table 6.** Performance of VOLE protocols (Arithmetic Circuit)

Protocol	Size	Speed	Non-Interactive
Wolverine	4	0.66 M	No
Mac'n'Cheese	3	0.4 M	No
QuickSilver	1	4.8 M	No
This work	$\frac{1}{b}$	$\geq 1.08$ M	Yes

There are other techniques for improving prover memory: Commit-and-prove to glue sub-circuits together (Lunar/Eclipse [15] [2]), streaming SNARKs (Gemini [11]). However, the reported constructions of these protocols require trusted setup (non-transparent), and the prover runtime cost of these protocols is magnitudes more expensive.

#### 6.4 Post-Quantum Alternative

It is possible to transform our protocol to a plausible post quantum design. In fact, there are only two modifications needed to make our protocol plausible post quantum. First, as we noted earlier, it is not a required to use Pedersen commitment to store committed inputs. It is entirely likely that we can find a post quantum commitment/encryption scheme and design a post-quantum input-mapping verification protocol to confirm the transformation. Second, it is possible to switch from a group based PCS scheme to a plausible post-quantum PCS scheme. We suspect the post-quantum version of our protocol offers comparable prover and verifier runtimes performances compared to that of the current version of our protocol. The communication cost, however, will likely be higher in the post-quantum setting in the default setting for “b” value, but flattens out when we move to smaller “b” value.

#### 6.5 Recent Developments

Very recently, major advances in multivariate polynomial commitment schemes has been made. One intriguing example is Binius, which is promising especially for proving binary field operations. However, it is not clear how much efficiency gain it can offer for integer arithmetics over u16/u32/u64 given the added complexity it may require to support it. Even if we count the cost of PCS schemes to zero for protocols leverages multivariate PCS schemes in exchange for more expensive communication cost, there are still 25% - 50% of prover cost it can't reduce and therefore cannot get even close to the benchmark performance our protocol offers.

Please refer to the supplemental material for details on the processing and tuning of boolean circuits.

### Acknowledgments

Special thanks to Danli Xie and Jam Jia for helping out on the coding and making the code more efficient.

### References

1. Ames, S., Hazay, C., Ishai, Y., Venkatasubramanian, M.: Liger: Lightweight sublinear arguments without a trusted setup. In: Thuraisingham, B.M., Evans, D., Malkin, T., Xu, D. (eds.) ACM CCS 2017. pp. 2087–2104. ACM Press, Dallas, TX, USA (Oct 31 – Nov 2, 2017). <https://doi.org/10.1145/3133956.3134104>



2. Aranha, D.F., Benedsen, E.M., Campanelli, M., Ganesh, C., Orlandi, C., Takahashi, A.: ECLIPSE: Enhanced compiling method for pedersen-committed zkSNARK engines. In: Hanaoka, G., Shikata, J., Watanabe, Y. (eds.) PKC 2022, Part I. LNCS, vol. 13177, pp. 584–614. Springer, Cham, Switzerland, Virtual Event (Mar 8–11, 2022). [https://doi.org/10.1007/978-3-030-97121-2\\_21](https://doi.org/10.1007/978-3-030-97121-2_21)
3. Arora, S., Lund, C., Motwani, R., Sudan, M., Szegedy, M.: Proof verification and hardness of approximation problems. In: 33rd FOCS. pp. 14–23. IEEE Computer Society Press, Pittsburgh, PA, USA (Oct 24–27, 1992). <https://doi.org/10.1109/SFCS.1992.267823>
4. Arora, S., Safra, S.: Probabilistic checking of proofs; A new characterization of NP. In: 33rd FOCS. pp. 2–13. IEEE Computer Society Press, Pittsburgh, PA, USA (Oct 24–27, 1992). <https://doi.org/10.1109/SFCS.1992.267824>
5. Babai, L., Fortnow, L., Levin, L.A., Szegedy, M.: Checking computations in polylogarithmic time. In: 23rd ACM STOC. pp. 21–31. ACM Press, New Orleans, LA, USA (May 6–8, 1991). <https://doi.org/10.1145/103418.103428>
6. Babai, L., Fortnow, L., Lund, C.: Non-deterministic exponential time has two-prover interactive protocols. In: 31st FOCS. pp. 16–25. IEEE Computer Society Press, St. Louis, MO, USA (Oct 22–24, 1990). <https://doi.org/10.1109/SFCS.1990.89520>
7. Baum, C., Malozemoff, A.J., Rosen, M.B., Scholl, P.: Mac’n’cheese: Zero-knowledge proofs for boolean and arithmetic circuits with nested disjunctions. In: Malkin, T., Peikert, C. (eds.) CRYPTO 2021, Part IV. LNCS, vol. 12828, pp. 92–122. Springer, Cham, Switzerland, Virtual Event (Aug 16–20, 2021). [https://doi.org/10.1007/978-3-030-84259-4\\_4](https://doi.org/10.1007/978-3-030-84259-4_4)
8. Ben-Sasson, E., Bentov, I., Horesh, Y., Riabzev, M.: Scalable zero knowledge with no trusted setup. In: Boldyreva, A., Micciancio, D. (eds.) CRYPTO 2019, Part III. LNCS, vol. 11694, pp. 701–732. Springer, Cham, Switzerland, Santa Barbara, CA, USA (Aug 18–22, 2019). [https://doi.org/10.1007/978-3-030-26954-8\\_23](https://doi.org/10.1007/978-3-030-26954-8_23)
9. Bhaduria, R., Fang, Z., Hazay, C., Venkatasubramanian, M., Xie, T., Zhang, Y.: Liger++: A new optimized sublinear IOP. In: Ligatti, J., Ou, X., Katz, J., Vigna, G. (eds.) ACM CCS 2020. pp. 2025–2038. ACM Press, Virtual Event, USA (Nov 9–13, 2020). <https://doi.org/10.1145/3372297.3417893>
10. Bootle, J., Cerulli, A., Chaidos, P., Groth, J., Petit, C.: Efficient zero-knowledge arguments for arithmetic circuits in the discrete log setting. In: Fischlin, M., Coron, J.S. (eds.) EUROCRYPT 2016, Part II. LNCS, vol. 9666, pp. 327–357. Springer, Berlin, Heidelberg, Germany, Vienna, Austria (May 8–12, 2016). [https://doi.org/10.1007/978-3-662-49896-5\\_12](https://doi.org/10.1007/978-3-662-49896-5_12)
11. Bootle, J., Chiesa, A., Hu, Y., Orrù, M.: Gemini: Elastic SNARKs for diverse environments. In: Dunkelman, O., Dziembowski, S. (eds.) EUROCRYPT 2022, Part II. LNCS, vol. 13276, pp. 427–457. Springer, Cham, Switzerland, Trondheim, Norway (May 30 – Jun 3, 2022). [https://doi.org/10.1007/978-3-031-07085-3\\_15](https://doi.org/10.1007/978-3-031-07085-3_15)
12. Boyle, E., Couteau, G., Gilboa, N., Ishai, Y.: Compressing vector OLE. In: Lie, D., Mannan, M., Backes, M., Wang, X. (eds.) ACM CCS 2018. pp. 896–912. ACM Press, Toronto, ON, Canada (Oct 15–19, 2018). <https://doi.org/10.1145/3243734.3243868>
13. Boyle, E., Couteau, G., Gilboa, N., Ishai, Y., Kohl, L., Rindal, P., Scholl, P.: Efficient two-round OT extension and silent non-interactive secure computation. In: Cavallaro, L., Kinder, J., Wang, X., Katz, J. (eds.) ACM CCS 2019. pp. 291–308. ACM Press, London, UK (Nov 11–15, 2019). <https://doi.org/10.1145/3319535.3354255>
14. Boyle, E., Couteau, G., Gilboa, N., Ishai, Y., Kohl, L., Scholl, P.: Efficient pseudorandom correlation generators: Silent OT extension and more. In: Boldyreva, A., Micciancio, D. (eds.) CRYPTO 2019, Part III. LNCS, vol. 11694, pp. 489–518. Springer, Cham, Switzerland, Santa Barbara, CA, USA (Aug 18–22, 2019). [https://doi.org/10.1007/978-3-030-26954-8\\_16](https://doi.org/10.1007/978-3-030-26954-8_16)
15. Campanelli, M., Faonio, A., Fiore, D., Querol, A., Rodríguez, H.: Lunar: A toolbox for more efficient universal and updatable zkSNARKs and commit-and-prove extensions. In: Tibouchi, M., Wang, H. (eds.) ASIACRYPT 2021, Part III. LNCS, vol. 13092, pp. 3–33. Springer, Cham, Switzerland, Singapore (Dec 6–10, 2021). [https://doi.org/10.1007/978-3-030-92078-4\\_1](https://doi.org/10.1007/978-3-030-92078-4_1)
16. Chen, B., Bünz, B., Boneh, D., Zhang, Z.: HyperPlonk: Plonk with linear-time prover and high-degree custom gates. In: Hazay, C., Stam, M. (eds.) EUROCRYPT 2023, Part II. LNCS, vol. 14005, pp. 499–530. Springer, Cham, Switzerland, Lyon, France (Apr 23–27, 2023). [https://doi.org/10.1007/978-3-031-30617-4\\_17](https://doi.org/10.1007/978-3-031-30617-4_17)
17. Chiesa, A., Hu, Y., Maller, M., Mishra, P., Vesely, P., Ward, N.P.: Marlin: Preprocessing zkSNARKs with universal and updatable SRS. In: Canteaut, A., Ishai, Y. (eds.) EUROCRYPT 2020, Part I.

- LNCS, vol. 12105, pp. 738–768. Springer, Cham, Switzerland, Zagreb, Croatia (May 10–14, 2020). [https://doi.org/10.1007/978-3-030-45721-1\\_26](https://doi.org/10.1007/978-3-030-45721-1_26)
18. Choudhuri, A.R., Garg, S., Goel, A., Sekar, S., Sinha, R.: SublonK: Sublinear prover PlonK. Cryptology ePrint Archive, Report 2023/902 (2023), <https://eprint.iacr.org/2023/902>
  19. Cramer, R., Damgård, I.: Zero-knowledge proofs for finite field arithmetic; or: Can zero-knowledge be for free? In: Krawczyk, H. (ed.) CRYPTO’98. LNCS, vol. 1462, pp. 424–441. Springer, Berlin, Heidelberg, Germany, Santa Barbara, CA, USA (Aug 23–27, 1998). <https://doi.org/10.1007/BFb0055745>
  20. Diamond, B.E., Posen, J.: Polylogarithmic proofs for multilinear towers. Cryptology ePrint Archive, Report 2024/504 (2024), <https://eprint.iacr.org/2024/504>
  21. Dittmer, S., Ishai, Y., Lu, S., Ostrovsky, R.: Improving line-point zero knowledge: Two multiplications for the price of one. In: Yin, H., Stavrou, A., Cremers, C., Shi, E. (eds.) ACM CCS 2022. pp. 829–841. ACM Press, Los Angeles, CA, USA (Nov 7–11, 2022). <https://doi.org/10.1145/3548606.3559385>
  22. Dittmer, S., Ishai, Y., Ostrovsky, R.: Line-point zero knowledge and its applications. Cryptology ePrint Archive, Report 2020/1446 (2020), <https://eprint.iacr.org/2020/1446>
  23. Frederiksen, T.K., Nielsen, J.B., Orlandi, C.: Privacy-free garbled circuits with applications to efficient zero-knowledge. In: Oswald, E., Fischlin, M. (eds.) EUROCRYPT 2015, Part II. LNCS, vol. 9057, pp. 191–219. Springer, Berlin, Heidelberg, Germany, Sofia, Bulgaria (Apr 26–30, 2015). [https://doi.org/10.1007/978-3-662-46803-6\\_7](https://doi.org/10.1007/978-3-662-46803-6_7)
  24. Gabizon, A., Williamson, Z.J., Ciobotaru, O.: PLONK: Permutations over lagrange-bases for oecumenical noninteractive arguments of knowledge. Cryptology ePrint Archive, Report 2019/953 (2019), <https://eprint.iacr.org/2019/953>
  25. Giacomelli, I., Madsen, J., Orlandi, C.: ZKBoo: Faster zero-knowledge for Boolean circuits. In: Holz, T., Savage, S. (eds.) USENIX Security 2016. pp. 1069–1083. USENIX Association, Austin, TX, USA (Aug 10–12, 2016)
  26. Goldwasser, S., Micali, S., Rackoff, C.: The knowledge complexity of interactive proof-systems (extended abstract). In: 17th ACM STOC. pp. 291–304. ACM Press, Providence, RI, USA (May 6–8, 1985). <https://doi.org/10.1145/22145.22178>
  27. Groth, J., Kohlweiss, M., Maller, M., Meiklejohn, S., Miers, I.: Updatable and universal common reference strings with applications to zk-SNARKs. In: Shacham, H., Boldyreva, A. (eds.) CRYPTO 2018, Part III. LNCS, vol. 10993, pp. 698–728. Springer, Cham, Switzerland, Santa Barbara, CA, USA (Aug 19–23, 2018). [https://doi.org/10.1007/978-3-319-96878-0\\_24](https://doi.org/10.1007/978-3-319-96878-0_24)
  28. Heath, D., Kolesnikov, V.: Stacked garbling for disjunctive zero-knowledge proofs. In: Canteaut, A., Ishai, Y. (eds.) EUROCRYPT 2020, Part III. LNCS, vol. 12107, pp. 569–598. Springer, Cham, Switzerland, Zagreb, Croatia (May 10–14, 2020). [https://doi.org/10.1007/978-3-030-45727-3\\_19](https://doi.org/10.1007/978-3-030-45727-3_19)
  29. Kiayias, A., Tang, Q.: How to keep a secret: leakage deterring public-key cryptosystems. In: Sadeghi, A.R., Gligor, V.D., Yung, M. (eds.) ACM CCS 2013. pp. 943–954. ACM Press, Berlin, Germany (Nov 4–8, 2013). <https://doi.org/10.1145/2508859.2516691>
  30. Maller, M., Bowe, S., Kohlweiss, M., Meiklejohn, S.: Sonic: Zero-knowledge SNARKs from linear-size universal and updatable structured reference strings. In: Cavallaro, L., Kinder, J., Wang, X., Katz, J. (eds.) ACM CCS 2019. pp. 2111–2128. ACM Press, London, UK (Nov 11–15, 2019). <https://doi.org/10.1145/3319535.3339817>
  31. Schnorr, C.P.: Efficient identification and signatures for smart cards. In: Brassard, G. (ed.) CRYPTO’89. LNCS, vol. 435, pp. 239–252. Springer, New York, USA, Santa Barbara, CA, USA (Aug 20–24, 1990). [https://doi.org/10.1007/0-387-34805-0\\_2](https://doi.org/10.1007/0-387-34805-0_2)
  32. Schoppmann, P., Gascón, A., Reichert, L., Raykova, M.: Distributed vector-OLE: Improved constructions and implementation. In: Cavallaro, L., Kinder, J., Wang, X., Katz, J. (eds.) ACM CCS 2019. pp. 1055–1072. ACM Press, London, UK (Nov 11–15, 2019). <https://doi.org/10.1145/3319535.3363228>
  33. Setty, S.: Spartan: Efficient and general-purpose zkSNARKs without trusted setup. In: Micciancio, D., Ristenpart, T. (eds.) CRYPTO 2020, Part III. LNCS, vol. 12172, pp. 704–737. Springer, Cham, Switzerland, Santa Barbara, CA, USA (Aug 17–21, 2020). [https://doi.org/10.1007/978-3-030-56877-1\\_25](https://doi.org/10.1007/978-3-030-56877-1_25)
  34. Setty, S., Lee, J.: Quarks: Quadruple-efficient transparent zkSNARKs. Cryptology ePrint Archive, Report 2020/1275 (2020), <https://eprint.iacr.org/2020/1275>
  35. Wahby, R.S., Tzialla, I., Shelat, A., Thaler, J., Walfish, M.: Doubly-efficient zkSNARKs without trusted setup. Cryptology ePrint Archive, Report 2017/1132 (2017), <https://eprint.iacr.org/2017/1132>
  36. Weng, C., Yang, K., Katz, J., Wang, X.: Wolverine: Fast, scalable, and communication-efficient zero-knowledge proofs for boolean and arithmetic circuits. In: 2021 IEEE Symposium on Security and Pri-

- vac. pp. 1074–1091. IEEE Computer Society Press, San Francisco, CA, USA (May 24–27, 2021). <https://doi.org/10.1109/SP40001.2021.00056>
37. Weng, C., Yang, K., Yang, Z., Xie, X., Wang, X.: AntMan: Interactive zero-knowledge proofs with sublinear communication. In: Yin, H., Stavrou, A., Cremers, C., Shi, E. (eds.) ACM CCS 2022. pp. 2901–2914. ACM Press, Los Angeles, CA, USA (Nov 7–11, 2022). <https://doi.org/10.1145/3548606.3560667>
  38. Yang, K., Sarkar, P., Weng, C., Wang, X.: QuickSilver: Efficient and affordable zero-knowledge proofs for circuits and polynomials over any field. In: Vigna, G., Shi, E. (eds.) ACM CCS 2021. pp. 2986–3001. ACM Press, Virtual Event, Republic of Korea (Nov 15–19, 2021). <https://doi.org/10.1145/3460120.3484556>
  39. Yang, K., Weng, C., Lan, X., Zhang, J., Wang, X.: Ferret: Fast extension for correlated OT with small communication. In: Ligatti, J., Ou, X., Katz, J., Vigna, G. (eds.) ACM CCS 2020. pp. 1607–1626. ACM Press, Virtual Event, USA (Nov 9–13, 2020). <https://doi.org/10.1145/3372297.3417276>
  40. Zhang, J., Xie, T., Zhang, Y., Song, D.: Transparent polynomial delegation and its applications to zero knowledge proof. In: 2020 IEEE Symposium on Security and Privacy. pp. 859–876. IEEE Computer Society Press, San Francisco, CA, USA (May 18–21, 2020). <https://doi.org/10.1109/SP40000.2020.00052>

## Appendix

### A. Univariate Polynomial Commitment Scheme

When a lot of inputs have value 0, it is possible that some  $z_i$  (except  $z_0$ , the sum of product of blinding keys) also computes value “0”, which may leak information about inputs.

The reason we didn’t pick the more popular FRI-based PCS scheme is because the polynomial our protocol needs to evaluate is expected to be small, and the prover and verifier runtime performance is automatically good. Instead, we use the PCS defined by Bootle et. al. which is depicted below with our terminology:

$$f(X) = (1 \ X^n \dots \ X^{(m-1)n}) \begin{pmatrix} z_{0,0} & z_{0,1} & \cdot & \cdot & z_{0,n-1} \\ z_{1,0} & z_{1,1} & \cdot & \cdot & z_{1,n-1} \\ \cdot & \cdot & & & \cdot \\ \cdot & \cdot & & & \cdot \\ \cdot & \cdot & & & \cdot \\ z_{m-1,0} & z_{m-1,1} & \cdot & \cdot & z_{m-1,n-1} \end{pmatrix} \begin{pmatrix} 1 \\ X \\ \cdot \\ \cdot \\ \cdot \\ X^{n-1} \end{pmatrix}$$

Figure 4

$\vec{z}$  are coefficients of the polynomial and are arranged to an  $m \times n$  matrix. The idea behind Bootle et al’s protocol is that the prover commits to the rows of this matrix using commitments  $C_0, \dots, C_{m-1}$  s.t.  $C_i = \sum_{j=0}^{n-1} z_{i,j} \cdot x^j$ . When given an evaluation point  $x$ , we use the homomorphic property of the commitment scheme to compute the commitment  $\prod_{i=0}^{m-1} C_i^{i \cdot n}$  to the vector:

$$\vec{c} = (1 \ x^n \dots \ x^{(m-1)n}) \begin{pmatrix} z_{0,0} & z_{0,1} & \cdot & \cdot & z_{0,n-1} \\ z_{1,0} & z_{1,1} & \cdot & \cdot & z_{1,n-1} \\ \cdot & \cdot & & & \cdot \\ \cdot & \cdot & & & \cdot \\ \cdot & \cdot & & & \cdot \\ z_{m-1,0} & z_{m-1,1} & \cdot & \cdot & z_{m-1,n-1} \end{pmatrix}$$

Figure 5

The prover opens this later commitment so that computing  $y = f(x)$  becomes trivial. To avoid leaking partial information about the coefficients of  $f(x)$  through these openings, Bootle et al.’s protocol inserts blinding values  $u_1, \dots, u_{n-1}$  to hide the weighted sum of the coefficients in each column

of the matrix, and make sure these blinding values cancel each other out when the polynomial gets evaluated at point  $x$ .

$$\begin{pmatrix} z_{0,0} & z_{0,1} + u_1 & \cdot & \cdot & z_{0,n-2} + u_{n-2} & z_{0,n-1} + u_{n-1} \\ z_{1,0} & z_{1,1} & \cdot & \cdot & z_{1,n-2} & z_{1,n-1} \\ \cdot & \cdot & & & \cdot & \cdot \\ \cdot & \cdot & & & \cdot & \cdot \\ \cdot & \cdot & & & \cdot & \cdot \\ z_{m-1,0} & z_{m-1,1} & \cdot & \cdot & z_{m-1,n-2} & z_{m-1,n-1} \\ -u_{q1} & -u_{q2} & \cdot & \cdot & -u_{qn-1} & 0 \end{pmatrix}$$

Figure 6

We do not need to make any modifications to the protocol itself. However, we need to redefine the domain of each variable in the matrix above since all coefficients are in  $\mathbb{Z}_q$  and cannot overflow  $p$  when computing commitments  $C_0, \dots, C_{m-1}$ . When given an evaluation point  $x$ , we use the homomorphic property of the commitment scheme to compute the commitment  $\prod_{i=0}^{m-1} C_i^{i \cdot n}$ .

Since coefficients  $z$  are in  $\mathbb{Z}_q$  and the exponentiation operations performed in group cannot overflow the order of  $\mathbb{G}$ , the verifier must compute powers of  $x$  (e.g.  $x, x^2, \dots, x^s$ ) in  $\mathbb{Z}_q$  before applying them to  $C_i$ .

Each blinding key  $u_j$  is randomly sampled from  $\mathbb{Z}_p$ , but we can limit each  $|u_j|$  to a  $|\sum_{i=0}^{m-1} z_{i,j} \cdot x^{j \cdot n}| + 80$ -bits value, which is approximately  $\approx 207$ -bits in our case for  $n = m = 2^5$ , enough to hide coefficients of  $T_i$ . For the first column  $j = 0$ , we expand  $z_0$  from a 61-bit value to an equivalent 207-bit value s.t.  $z_0 = z_0 + s \cdot q$  for some 146-bit random  $s$ .

$u_j$  will get cancelled by  $u_{qj} = u_j \bmod q$  once the domain of the final evaluation is reduced to domain  $\mathbb{Z}_q$ .

## B. The sub-protocol for boolean circuit validation

A common requirement in proving boolean circuits is to enforce input data  $b_i \in \{0, 1\}$  for  $i \in \{1, \dots, l\}$  (this is not new, but we need this sub-protocol defined to make our main protocol easier to parse), such relation is defined as:

$$\{(\vec{b}', \in \mathbb{Z}_q^l; \vec{b}, \vec{\beta} \in \mathbb{Z}_q^l) : b'_i = b_i + X\beta_i \wedge \forall b_i \in [0, 1] \wedge \forall_i \in [1, \dots, l]\} \quad (76)$$

In practice, it is useful to decompose  $l$  full integer inputs into  $l \cdot 32$  bits (assuming we use 32 bits to represent a full integer). If a committed value  $b_i$  is in  $[0, 1]$ , then its linear polynomial form  $b_i$  must have the following property:

$$(b'_i \cdot b'_i - b'_i) = \delta_{1,i}x + \delta_{2,i}x^2 \quad (77)$$

Where  $\delta_{2,i} = \beta_i^2$ , and  $\delta_{1,i} = \beta_i$  when  $b_i$  is 1 and  $\delta_{1,i} = -\beta$  when  $b_i$  is 0. To prove the correctness for all  $b_i \in \{0, 1\}$ , the prover commits to polynomials  $D_1, D_2$ :

$$D_1 = u_1^{\delta_{1,1}} u_2^{\delta_{1,2}} \dots u_l^{\delta_{1,l}} h^{\rho_1}, \quad D_2 = u_1^{\delta_{2,1}} u_2^{\delta_{2,2}} \dots u_l^{\delta_{2,l}} h^{\rho_2}$$

Where  $D_1$  commits to coefficients on the  $x$  term and  $D_2$  commits to coefficients on the  $x^2$  term. The prover can easily join two polynomial commitments into one and sends only one element  $D$  to the verifier.

$$D = \prod_{i=1}^l u_i^{\delta_{1i}} \cdot \prod_{i=1}^l u_{i+l}^{\delta_{2i}} \cdot h^\rho \in \mathbb{G} \quad (78)$$

When the challenge  $k$  is received, the prover sends the evaluation results  $y_1, y_2$  to the verifier, and then engages with the verifier to verify the correctness of  $y_1, y_2$  at point  $k$ , and checks if the equality below is true:

$$y_1 \cdot X + y_2 \cdot X^2 = \sum_{j=1}^{l \cdot 32} (b'_i \cdot b'_i - b'_i) \cdot k^i \quad (79)$$

Once we know all linear polynomials map to either 0 or 1, it is trivial to recombine the linear polynomial form of a full integer input  $a'_i$  from 32 decomposed bits  $b'_{i,j}$  for  $j = \{1, \dots, 32\}$ .

$$a'_i = \sum_{j=1}^{32} b'_{i,j} \cdot 2^j \quad (80)$$

We define the protocol BooleanityTest using two sub-protocols:

*Input* :  $(\vec{b}' \in \mathbb{Z}_q; \vec{b}, \vec{\beta}' \in \mathbb{Z}_q)$

*P compute* :

$l = |\vec{b}'| \in \mathbb{Z}_q$

$\rho \xleftarrow{\$} \mathbb{Z}_q$

$\delta_{1i} = b_i \beta_i + b_i \beta_i - \beta_i \in \mathbb{Z}_q \quad i = \{1, \dots, l\}$

$\delta_{2i} = \beta_i^2 \in \mathbb{Z}_q \quad i = \{1, \dots, l\}$

$D = \prod_{i=1}^l u_i^{\delta_{1,i}} \cdot \prod_{i=1}^l u_{i+l}^{\delta_{2,i}} \cdot h^\rho \in \mathbb{G}$

Protocol BooleanityTest-Setup

After a challenge  $x$  is sent from the verifier to the prover, the protocol moves to the verification stage defined by the following sub-protocol.

*Input* :  $(\vec{b}' \in \mathbb{Z}_q; \rho, \vec{\delta}_1, \vec{\delta}_2' \in \mathbb{Z}_q)$

*P's input* :  $(\vec{b}'; \rho, \vec{\delta}_1, \vec{\delta}_2')$ ; *V's input* :  $(\vec{b}')$

$\mathcal{V} \rightarrow \mathcal{P}$  :  $k \xleftarrow{\$} \mathbb{Z}_q$

*P compute* :

$y_1 = \sum_{i=1}^l \delta_{1,i} k^i \in \mathbb{Z}_q, \quad y_2 = \sum_{i=1}^l \delta_{2,i} k^i \in \mathbb{Z}_q$

$\mathcal{P} \rightarrow \mathcal{V}$  :  $y_1, y_2$

$\mathcal{P}, \mathcal{V}$  engage to evaluate :

**if**  $PolyEval(D, y_1 + y_2 k^l, k; \vec{\delta}_1 || \vec{\delta}_2, \rho)$

$\wedge y_1 \cdot x + y_2 \cdot x^2 \stackrel{?}{=} \sum_{j=1}^{l \cdot 32} (b'_i \cdot b'_i - b'_i) \cdot k^i$

**return true**

**else return false**

## Protocol BooleanityTest-Verify

*Proof.* Perfect completeness follows from the fact that the protocol BooleanityTest is trivially complete.

To prove PHVZK for relation 76, we define a simulator  $\mathcal{S}_{b-test}$ . To start,  $\mathcal{S}_{b-test}$  randomly generates a group element  $D$ , which represents the committed polynomial. After challenge  $x$  is received from the verifier,

$\mathcal{S}_{b-test}$  uses a simulator  $\mathcal{S}_p$  to simulate proof transcripts needed for polynomial commitment evaluation, which we know exist for a fact [10].

The simulators  $\mathcal{S}_{b-test}$  and  $\mathcal{S}_p$  choose all proof elements and challenges according to the randomness supplied by the adversary from their respective domains or compute them directly as described in the protocol. Since all elements in proof transcripts are either independently randomly distributed or their relationship is fully defined by the verification equations, we can conclude that protocol BooleanityTest is PHVZK.

To prove this protocol has CWEE, we first define an extractor  $\mathcal{X}_{b-test}$  for Protocol Booleanity that extracts witnesses  $\vec{b}, \vec{\beta}$ .

The extractor first receives a vector commitment  $D$  from the prover. The extractor then generates  $2l + 1$  challenges  $k_1, \dots, k_{2l+1}$  and retrieves  $2l + 1$   $\vec{y}_1$  and  $2l + 1$   $\vec{y}_2$  through repeated rewinding. The extractor  $\mathcal{X}_{b-test}$  then calls on an extractor for the polynomial commitment evaluation protocol (which we know exists for a fact) to extract witnesses  $\vec{\delta}_1, \vec{\delta}_2, \phi$  or else we find a non-trivial relationship between elements in  $\vec{u}, h$ .

It is trivial to extract witnesses  $b_i, \beta_i$  from  $\delta_1, \delta_2$  for  $i = \{1, \dots, l\}$  s.t. they must pass the equality test defined in equation 79 except with negligible probability of a dishonest prover making the right guess on  $x$