

MATRIX ADDITION, HADAMARD PRODUCTS, AND MATRIX MULTIPLICATION

1	INTRODUCTION	1
2	MATRIX ADDITION: THE PROCESS	1
2.1	Steps in the process.	2
2.2	Commentary on each step.	2
3	HADAMARD PRODUCT: THE PROCESS	3
3.1	Steps in the process.	4
3.2	Commentary on each step.	4
4	MATRIX MULTIPLICATION: THE PROCESS	5
4.1	Steps in the process.	5
4.2	Commentary on each step.	6
5	THE MATHS	8
6	THE CODE	8
7	EXAMPLES	8
	REFERENCES	9

1. INTRODUCTION

Matrix multiplication is a core operation in the forward pass of neural networks, underlying the computations performed by fully connected layers—where input activations are combined with learned weights to produce subsequent activations. The Hadamard product—defined as the elementwise product of two matrices of the same shape—also plays an essential role, particularly in gating mechanisms common in recurrent networks and attention modules, where it modulates information flow on a per-coordinate basis.

In this document, we develop methods for verifying integer relations involving these fundamental operations and their linear combinations within arithmetic circuits. Specifically, we consider integer-valued matrices without assuming any scaling or quantization. (Verification involving matrices scaled for fixed-point representations—known as *quantized matrix multiplication*—is treated in a separate companion document [1].)

The overarching verification strategy employed throughout this document is as follows: we begin with an integer equality that we aim to verify within an arithmetic circuit over a finite field $\mathbb{Z}/p\mathbb{Z}$. Directly enforcing integer equalities is impossible in modular arithmetic, as a congruence modulo p only ensures equality up to some multiple of p . Therefore, the core verification challenge is to guarantee—through the application of carefully chosen inequalities and range checks—that this multiple of p must be zero. Since all integers are first converted into least-residue representatives modulo p *before* entering the circuit, the circuit itself cannot distinguish between distinct integers differing by multiples of p . Consequently, establishing that inputs lie within a predetermined, known set of residues modulo p necessarily relies on off-circuit justifications and assumptions. This interplay between modular constraints (enforced within the circuit) and integer range assumptions (justified off-circuit) is fundamental to all arithmetic circuit verifications.

With this verification paradigm in place, we begin by examining the simplest scenario—*matrix addition*—which establishes the basic structure and reasoning behind our verification process. Subsequently, we consider the *Hadamard product*, structurally analogous to addition but involving elementwise multiplication. Finally, we address *matrix multiplication* in its “vanilla” (non-quantized) form, as well as generalized expressions of the form

$$\alpha AB + \beta C,$$

where A , B , and C are integer matrices of compatible dimensions, and α, β are known integer scalars.

Together, these foundational operations encompass a wide variety of neural network architectures—from simple affine transformations to sophisticated gating and attention mechanisms—enabling efficient and rigorous zero-knowledge verification of neural network behavior.

2. MATRIX ADDITION: THE PROCESS

Let A_1, \dots, A_k be matrices of the same dimension $m \times n$, and let $\alpha_1, \dots, \alpha_k$ be integer scalars. We define the weighted sum

$$B = \alpha_1 A_1 + \dots + \alpha_k A_k,$$

with $B = [b_{ij}]$, $A_t = [a_{ij}^{(t)}]$ for $t = 1, \dots, k$, and

$$b_{ij} = \alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)}. \quad (2.1)$$

Verifying such a matrix addition within an arithmetic circuit amounts to checking mn independent linear constraints over $\mathbb{Z}/p\mathbb{Z}$. This operation arises naturally in affine layers, residual blocks, and many preprocessing or postprocessing steps in neural networks. The constraints are simple to implement: we verify that

$$b_{ij} \equiv \alpha_1 a_{ij}^{(1)} + \cdots + \alpha_k a_{ij}^{(k)} \pmod{p}.$$

However, we must ensure that the integer representation of each sum lies within a known range to interpret the result unambiguously.

2.1. Steps in the process.

- (1) The circuit operates over the finite field $\mathbb{Z}/p\mathbb{Z}$, where p is a prime.
- (2) Let $A_1 = [a_{ij}^{(1)}], \dots, A_k = [a_{ij}^{(k)}], B = [b_{ij}]$ be integer matrices of dimensions $m \times n$, and let $\alpha_1, \dots, \alpha_k$ be integer scalars. Assume there is an integer h such that for all i, j ,

$$h - p \leq \alpha_1 a_{ij}^{(1)} + \cdots + \alpha_k a_{ij}^{(k)}, b_{ij} < h. \quad (2.2)$$

- (3) The goal of the circuit is to verify that $\alpha_1 A_1 + \cdots + \alpha_k A_k = B$, that is, for all i, j ,

$$b_{ij} = \alpha_1 a_{ij}^{(1)} + \cdots + \alpha_k a_{ij}^{(k)}. \quad (2.3)$$

- (4) Most frameworks work exclusively with least residue representations. Accordingly, we use \bar{x} to denote the least residue of $x \pmod{p}$: for all i, j, ℓ ,

$$\bar{\alpha}_\ell \equiv \alpha_\ell \pmod{p}, \quad \bar{a}_{ij}^{(\ell)} \equiv a_{ij}^{(\ell)} \pmod{p}, \quad \bar{b}_{ij} \equiv b_{ij} \pmod{p}, \quad 0 \leq \bar{\alpha}_\ell, \bar{a}_{ij}^{(\ell)}, \bar{b}_{ij} < p. \quad (2.4)$$

- (5) For all i, j , impose the constraint

$$(\bar{\alpha}_1 \bar{a}_{ij}^{(1)} + \cdots + \bar{\alpha}_k \bar{a}_{ij}^{(k)}) - \bar{b}_{ij} \equiv 0 \pmod{p}. \quad (2.5)$$

- (6) Given that the bounds in (2.2) hold and the constraints (2.5) are satisfied for all i, j , we may conclude that the original integer relation (2.3) holds entrywise. That is, under our assumptions, the congruences (2.5) correctly enforce the desired matrix sum

$$\alpha_1 A_1 + \cdots + \alpha_k A_k = B. \quad (2.6)$$

2.2. Commentary on each step.

- (1) Typically, p is an n -bit prime satisfying

$$2^{n-1} \leq p < 2^n,$$

with $n \approx 256$. In practice, we use the prime field associated with the scalar field of the BN254 elliptic curve, a 254-bit prime that offers a good balance between security and efficiency. This field is widely supported in cryptographic applications and zk-SNARK frameworks due to the curve's pairing-friendly properties and efficient arithmetic over $\mathbb{Z}/p\mathbb{Z}$.

- (2) The offset parameter h defines an integer interval of length p within which each weighted sum and each b_{ij} is assumed to lie. Two common conventions are:

- the *least residue range* $[0, p)$, corresponding to $h = p$;
- the *balanced residue range* $(-p/2, p/2]$, corresponding to $h = (p + 1)/2$.

As with any field-based circuit, the circuit cannot distinguish between integers that differ by a multiple of p , so the bound

$$h - p \leq b_{ij} < h$$

must be justified off-circuit.

The same applies to the weighted sum. In many cases, we may assume the range bound

$$h - p \leq \alpha_1 a_{ij}^{(1)} + \cdots + \alpha_k a_{ij}^{(k)} < h$$

holds automatically due to the size of p and small magnitudes of the scalars and matrix entries. However, if explicit justification is needed, we may enforce a bound in-circuit by verifying that

$$a_{ij}^{(\ell)} \in [-U, U] \quad (1 \leq \ell \leq k),$$

for some positive integer U satisfying

$$(|\alpha_1| + \dots + |\alpha_k|)U < p/2.$$

This ensures that the sum lies in $(-p/2, p/2)$. However, as before, this reasoning assumes each $a_{ij}^{(\ell)}$ is represented within a known residue interval of length p . Without that, even a valid range check cannot be interpreted unambiguously. For further discussion of this subtlety and the implementation of range checks, see our companion documents [2, 3].

- (3) In other words, we wish to verify the correctness of mn independent weighted sums.
- (4) In typical usage, the prover supplies the least residues $\bar{A}_1, \dots, \bar{A}_k$, and \bar{B} as part of the witness, where each entry is the canonical representative modulo p . If any matrix \bar{A}_ℓ is fixed—e.g., a constant bias matrix—it may be incorporated into the circuit directly. Likewise, the scalars $\bar{\alpha}_\ell$ are the least residues of known integers α_ℓ and may be either public parameters or hardcoded constants, depending on the application.
- (5) Although $\bar{\alpha}_\ell \equiv \alpha_\ell \pmod{p}$, $\bar{a}_{ij}^{(\ell)} \equiv a_{ij}^{(\ell)} \pmod{p}$, and $\bar{b}_{ij} \equiv b_{ij} \pmod{p}$, we state the constraint in terms of the least residues to reflect what is literally enforced in the circuit. The prover supplies $\bar{a}_{ij}^{(\ell)}$ and \bar{b}_{ij} as field elements in $[0, p)$, and the constraint

$$\left(\bar{\alpha}_1 \bar{a}_{ij}^{(1)} + \dots + \bar{\alpha}_k \bar{a}_{ij}^{(k)} \right) - \bar{b}_{ij} \equiv 0 \pmod{p}$$

is enforced using native field operations in the arithmetic circuit. The distinction between x and \bar{x} is semantically irrelevant to the field, but helps clarify the origin of each quantity.

- (6) From (2.5) and (2.4), we have

$$\left(\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} \right) - b_{ij} \equiv 0 \pmod{p},$$

so there exists an integer t such that

$$\left(\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} \right) - b_{ij} = tp.$$

Using (2.2), we know

$$h - p \leq \left(\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} \right), b_{ij} \leq h - 1.$$

(All terms involved are integers, so a strict inequality like $b_{ij} < h$ is equivalent to $b_{ij} \leq h - 1$.) Subtracting gives

$$(h - p) - (h - 1) \leq \left(\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} \right) - b_{ij} \leq (h - 1) - (h - p),$$

i.e.,

$$-(p - 1) \leq \left(\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} \right) - b_{ij} \leq p - 1.$$

Hence $tp \in (-p, p)$, implying $t \in (-1, 1)$, so $t = 0$, and therefore

$$\alpha_1 a_{ij}^{(1)} + \dots + \alpha_k a_{ij}^{(k)} = b_{ij}.$$

3. HADAMARD PRODUCT: THE PROCESS

Let $A = [a_{ij}]$ and $B = [b_{ij}]$ be matrices of the same dimension $m \times n$. Their *Hadamard product*, denoted $A \odot B$, is the matrix $C = [c_{ij}]$, also of dimension $m \times n$, defined by

$$c_{ij} = a_{ij} b_{ij}. \tag{3.1}$$

This operation arises naturally in gating mechanisms and feature-wise modulations in neural networks, and is efficient to verify in constraint systems due to its lack of cross-term interactions.

Verifying a Hadamard product computation within an arithmetic circuit amounts to checking the validity of mn independent multiplication constraints. These constraints are straightforward to implement: for each i, j , we simply verify that

$$c_{ij} \equiv a_{ij} b_{ij} \pmod{p}$$

However, careful attention must be paid to the range and encoding of the input and output values to ensure correctness at the integer level.

We consider a slightly more general setting: given integer matrices A, B, C , and D of equal dimensions and integer scalars α, β , we verify that

$$\alpha(A \odot B) + \beta C = D.$$

This generalization incurs little additional complexity in the constraint system but captures a broader class of common structures, including residual connections and affine gates.

3.1. Steps in the process. Each step in the process has a corresponding explanatory note in Subsection 3.2 that provides additional context and details.

- (1) The circuit operates over the finite field $\mathbb{Z}/p\mathbb{Z}$, where p is a prime.
- (2) Let $A = [a_{ij}], B = [b_{ij}], C = [c_{ij}], D = [d_{ij}]$ be integer matrices of dimensions $m \times n$. Let α and β be integers. Assume there is an integer h such that for all i, j ,

$$h - p \leq \alpha a_{ij} b_{ij} + \beta c_{ij}, d_{ij} < h. \quad (3.2)$$

- (3) The goal of the circuit is to verify that $\alpha(A \odot B) + \beta C = D$. That is, for all i, j ,

$$\alpha a_{ij} b_{ij} + \beta c_{ij} = d_{ij}. \quad (3.3)$$

- (4) Most frameworks work exclusively with least residue representations. Accordingly, we use \bar{x} to denote the least residue of $x \bmod p$:

$$\bar{\alpha} \equiv \alpha \bmod p, \quad \bar{\beta} \equiv \beta \bmod p, \quad 0 \leq \bar{\alpha}, \bar{\beta} < p,$$

and for all i, j ,

$$\bar{a}_{ij} \equiv a_{ij} \bmod p, \quad \bar{b}_{ij} \equiv b_{ij} \bmod p, \quad \bar{c}_{ij} \equiv c_{ij} \bmod p, \quad \bar{d}_{ij} \equiv d_{ij} \bmod p, \quad 0 \leq \bar{a}_{ij}, \bar{b}_{ij}, \bar{c}_{ij}, \bar{d}_{ij} < p. \quad (3.4)$$

- (5) For all i, j , impose the constraint

$$\bar{\alpha} \bar{a}_{ij} \bar{b}_{ij} + \bar{\beta} \bar{c}_{ij} - \bar{d}_{ij} \equiv 0 \bmod p. \quad (3.5)$$

- (6) Given that the bounds in (3.2) hold and the constraints (3.5) are satisfied for all i, j , we may conclude that the original integer relation (3.3) holds entrywise. That is, under our assumptions, the congruences (3.5) correctly enforce the desired Hadamard product

$$\alpha(A \odot B) + \beta C = D. \quad (3.6)$$

3.2. Commentary on each step.

- (1) See Comment (1) of Subsection 2.2.
- (2) The offset parameter h provides flexibility in defining an integer range of length p to which all values are assumed to belong. Prior to entering the circuit, we typically assume that each term $\alpha a_{ij} b_{ij} + \beta c_{ij}$ and each value d_{ij} lies within such a range. Two common conventions are:
 - the *least residue range* $[0, p)$, corresponding to $h = p$;
 - the *balanced residue range* $(-p/2, p/2]$, corresponding to $h = (p + 1)/2$.

Because arithmetic circuits over $\mathbb{Z}/p\mathbb{Z}$ cannot distinguish between integers that differ by a multiple of p , the condition

$$h - p \leq d_{ij} < h$$

cannot be enforced within the circuit itself. It must instead be justified through off-circuit reasoning, either by construction of the inputs or by assumption.

In most practical scenarios, the corresponding bound on the product,

$$h - p \leq \alpha a_{ij} b_{ij} + \beta c_{ij} < h,$$

can likely also be justified off-circuit. However, if such justification is not available, we may instead enforce the product bound in-circuit by performing a range check to verify that

$$a_{ij}, b_{ij}, c_{ij} \in [-U, U)$$

for some positive integer U satisfying $|\alpha|U^2 + |\beta|U < p/2$. In that case, it follows that

$$\alpha a_{ij} b_{ij} + \beta c_{ij} \in [-(|\alpha|U^2 + |\beta|U), |\alpha|U^2 + |\beta|U] \subseteq (-p/2, p/2),$$

so the product lies within the balanced residue range.

Still, the soundness of this range check implicitly relies on the assumption that a_{ij}, b_{ij}, c_{ij} are themselves represented in the balanced residue range $(-p/2, p/2]$, or some other known range of length p , to begin with. Without such an assumption, the result of the range check cannot be interpreted unambiguously. For further discussion of this subtlety and the implementation of range checks, see our companion documents [2, 3].

- (3) In other words, the goal is to verify mn independent multiplications.
- (4) In typical usage, the prover supplies the least residues $\bar{A} = [\bar{a}_{ij}]$, $\bar{B} = [\bar{b}_{ij}]$, $\bar{C} = [\bar{c}_{ij}]$, and $\bar{D} = [\bar{d}_{ij}]$ as part of the witness, where each entry is the canonical representative of its corresponding integer modulo p . Although it is possible to compute \bar{D} in-circuit from $\bar{A}, \bar{B}, \bar{C}$, doing so increases circuit size and is unnecessary for most modular zkML designs.

If B or C represent fixed matrices—such as model weights or bias terms—they may be incorporated directly into the circuit as constant values. This eliminates the need for the prover to supply \bar{B} or \bar{C} as part of the witness and can improve both prover performance and circuit compactness.

Likewise, the scalars $\bar{\alpha}$ and $\bar{\beta}$ may either be treated as public parameters or compiled into the circuit as fixed constants. In many practical scenarios—such as residual connections or affine gates—these coefficients are known in advance and can be hardcoded to minimize input complexity.

- (5) Although $\bar{\alpha} \equiv \alpha \pmod{p}$, and similarly for $\bar{\beta}, \bar{a}_{ij}, \bar{b}_{ij}, \bar{c}_{ij}$, and \bar{d}_{ij} , we state the constraint in terms of the least residue representatives to reflect what is literally enforced in the circuit. The prover supplies $\bar{a}_{ij}, \bar{b}_{ij}, \bar{c}_{ij}, \bar{d}_{ij} \in [0, p)$ as field elements, and the constraint $\bar{\alpha}\bar{a}_{ij}\bar{b}_{ij} + \bar{\beta}\bar{c}_{ij} - \bar{d}_{ij} \equiv 0 \pmod{p}$ is imposed using arithmetic over $\mathbb{Z}/p\mathbb{Z}$. The distinction between x and \bar{x} is not meaningful within the field, but it helps clarify the provenance of the values and emphasizes that no signed interpretation occurs within the circuit itself.
- (6) For all i, j , we have from (3.4) and (3.5) that

$$\alpha a_{ij} b_{ij} + \beta c_{ij} - d_{ij} \equiv 0 \pmod{p},$$

and hence there exists an integer t such that

$$\alpha a_{ij} b_{ij} + \beta c_{ij} - d_{ij} = tp. \quad (3.7)$$

By the assumption (3.2), we know that

$$h - p \leq \alpha a_{ij} b_{ij} + \beta c_{ij}, d_{ij} \leq h - 1.$$

(All terms involved are integers, so a strict inequality like $d_{ij} < h$ is equivalent to $d_{ij} \leq h - 1$.) Subtracting these bounds gives

$$(h - p) - (h - 1) \leq (\alpha a_{ij} b_{ij} + \beta c_{ij}) - d_{ij} \leq (h - 1) - (h - p),$$

which simplifies to

$$-(p - 1) \leq (\alpha a_{ij} b_{ij} + \beta c_{ij}) - d_{ij} \leq p - 1. \quad (3.8)$$

Substituting from (3.7), we conclude that $tp \in (-p, p)$, so $t \in (-1, 1)$, and hence $t = 0$. It follows that

$$\alpha a_{ij} b_{ij} + \beta c_{ij} = d_{ij}.$$

4. MATRIX MULTIPLICATION: THE PROCESS

Matrix multiplication lies at the heart of neural network computation. It is the fundamental operation performed by fully connected layers, where the output activations are computed as a linear transformation of the inputs using a learned weight matrix. Given integer matrices $A = [a_{ik}]$ of dimensions $\ell \times m$ and $B = [b_{kj}]$ of dimensions $m \times n$, their product is the matrix $C = [c_{ij}]$ of dimensions $\ell \times n$, with entries defined by

$$c_{ij} = \sum_{k=1}^m a_{ik} b_{kj}. \quad (4.1)$$

To verify this relation within an arithmetic circuit, we interpret the sum in (4.1) as a composition of multiplication and addition gates over a finite field. At a high level, we impose a constraint for each output entry c_{ij} that asserts it is equal to the inner product of the i -th row of A with the j -th column of B . As with the other operations in this document, we work over the finite field $\mathbb{Z}/p\mathbb{Z}$, and therefore must ensure that all inputs and outputs are appropriately encoded, and that the integer-level equality is soundly recovered from modular constraints using known bounds on the input values.

4.1. Steps in the process. Each step in the process has a corresponding explanatory note in Subsection 4.2 that provides additional context and details.

- (1) The circuit operates over the finite field $\mathbb{Z}/p\mathbb{Z}$, where p is a prime.
- (2) Let $A = [a_{ik}], B = [b_{kj}], C = [c_{ij}], D = [d_{ij}]$ be integer matrices of dimensions $\ell \times m, m \times n, \ell \times n$, respectively. Let α and β be integers.

(2a) Assume there is an integer h such that for all i, k, j ,

$$h - p \leq \alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij}, d_{ij} < h. \quad (4.2)$$

(2b) If this assumption cannot be justified, then assume there is an integer h such that for all i, k, j ,

$$h - p \leq a_{ik}, b_{kj}, c_{ij}, d_{ij} < h, \quad (4.3)$$

and perform a range check to ensure that

$$a_{ik}, b_{kj}, c_{ij} \in [-U, U), \quad (4.4)$$

for some positive integer U satisfying

$$|\alpha| m U^2 + |\beta| U < \min\{-(h - p), p\}. \quad (4.5)$$

(3) The goal of the circuit is to verify that $\alpha(AB) + \beta C = D$. That is, for all i, j ,

$$\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} = d_{ij}. \quad (4.6)$$

(4) Most frameworks work exclusively with least residue representations. Accordingly, we use \bar{x} to denote the least residue of $x \bmod p$:

$$\bar{\alpha} \equiv \alpha \bmod p, \quad \bar{\beta} \equiv \beta \bmod p, \quad 0 \leq \bar{\alpha}, \bar{\beta} < p,$$

and for all i, k, j ,

$$\bar{a}_{ik} \equiv a_{ik} \bmod p, \quad \bar{b}_{kj} \equiv b_{kj} \bmod p, \quad \bar{c}_{ij} \equiv c_{ij} \bmod p, \quad \bar{d}_{ij} \equiv d_{ij} \bmod p, \quad 0 \leq \bar{a}_{ik}, \bar{b}_{kj}, \bar{c}_{ij}, \bar{d}_{ij} < p. \quad (4.7)$$

(5) For all i, j , impose the constraint

$$\bar{\alpha} \left(\sum_{k=1}^m \bar{a}_{ik} \bar{b}_{kj} \right) + \bar{\beta} \bar{c}_{ij} - \bar{d}_{ij} \equiv 0 \bmod p. \quad (4.8)$$

(6) Given that the bounds in (4.2) hold and the constraints (4.8) are satisfied for all i, j , we may conclude that the original integer relation (4.6) holds entrywise. That is, under our assumptions, the congruences (3.5) correctly enforce the desired generalized matrix product

$$\alpha(AB) + \beta C = D. \quad (4.9)$$

4.2. Commentary on each step.

(1) See Comment (1) of Subsection 2.2.

(2) The offset parameter h provides flexibility in defining an integer interval of length p to which all values are assumed to belong. Prior to entering the circuit, we typically assume that each term $\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij}$ and each value d_{ij} lies within such a range. Two common conventions are:

- the *least residue range* $[0, p)$, corresponding to $h = p$;
- the *balanced residue range* $(-p/2, p/2]$, corresponding to $h = (p + 1)/2$.

Because arithmetic circuits over $\mathbb{Z}/p\mathbb{Z}$ cannot distinguish between integers that differ by a multiple of p , the condition

$$h - p \leq d_{ij} < h$$

must be justified off-circuit.

The same applies to the term $\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij}$. In many scenarios, such bounds may be justified off-circuit, particularly when m is small and the matrix entries are known to be bounded. However, when such justification is unavailable, we must instead enforce the bound in-circuit by verifying that

$$a_{ik}, b_{kj}, c_{ij} \in [-U, U)$$

for some positive integer U satisfying

$$|\alpha| m U^2 + |\beta| U < \min\{-(h-p), p\}.$$

This guarantees that

$$\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} \in [-(|\alpha| m U^2 + |\beta| U), |\alpha| m U^2 + |\beta| U] \subseteq (h-p, h). \quad (4.10)$$

As always, the validity of the range check depends on the off-circuit assumption that a_{ik} , b_{kj} , and c_{ij} are encoded using a known, consistent set of representatives modulo p (e.g., the balanced residue range). Without such an assumption, even a successful range check does not determine the underlying integer values uniquely. For further discussion of this subtlety and the implementation of range checks, see our companion documents [2, 3].

- (3) The goal is to verify ℓn constraints, one for each entry in the output matrix D . Each constraint corresponds to a generalized inner product, involving the weighted sum of a row of A and a column of B , followed by an affine shift by βc_{ij} . The structure resembles an inner product but includes both a scalar multiplier α and an additive offset.
- (4) In typical usage, the prover supplies the least residues $\bar{A} = [\bar{a}_{ik}]$, $\bar{B} = [\bar{b}_{kj}]$, $\bar{C} = [\bar{c}_{ij}]$, and $\bar{D} = [\bar{d}_{ij}]$ as part of the witness, where each entry is the canonical representative of its corresponding integer modulo p . Although it is possible to compute \bar{D} in-circuit from $\bar{A}, \bar{B}, \bar{C}$, doing so increases circuit size and is unnecessary in most modular zkML designs.

If B or C represent fixed matrices—such as learned model weights or constant bias terms—they may be hardcoded into the circuit. This eliminates the need for the prover to supply \bar{B} or \bar{C} and reduces prover-side workload.

Likewise, the scalars $\bar{\alpha}$ and $\bar{\beta}$ may be compiled into the circuit as fixed constants or supplied as public parameters, depending on the use case. In most inference settings, these coefficients are known in advance and are typically hardcoded to reduce input complexity and prover effort.

- (5) Although $\bar{\alpha} \equiv \alpha \pmod{p}$, and similarly for $\bar{\beta}, \bar{a}_{ik}, \bar{b}_{kj}, \bar{c}_{ij}$, and \bar{d}_{ij} , we express the constraint in terms of least residues to reflect the actual field-level expressions used in circuit implementation. The prover supplies $\bar{a}_{ij}, \bar{b}_{ij}, \bar{c}_{ij}, \bar{d}_{ij} \in [0, p)$ as field elements, and the constraint

$$\bar{\alpha} \left(\sum_{k=1}^m \bar{a}_{ik} \bar{b}_{kj} \right) + \bar{\beta} \bar{c}_{ij} - \bar{d}_{ij} \equiv 0 \pmod{p}$$

is enforced using arithmetic over $\mathbb{Z}/p\mathbb{Z}$. The distinction between x and \bar{x} is invisible to the circuit but clarifies how values are encoded at preprocessing time.

- (6) For all i, j , we have from (4.7) and (4.8) that

$$\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} - d_{ij} \equiv 0 \pmod{p},$$

and hence there exists an integer t such that

$$\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} - d_{ij} = tp. \quad (4.11)$$

From assumption (4.2) (alternatively, (4.10)), we know that

$$h-p \leq \alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij}, d_{ij} \leq h-1.$$

(All terms involved are integers, so a strict inequality like $d_{ij} < h$ is equivalent to $d_{ij} \leq h-1$.) Subtracting gives

$$(h-p) - (h-1) \leq \left(\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} \right) - d_{ij} \leq (h-1) - (h-p),$$

i.e.,

$$-(p-1) \leq \left(\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} \right) - d_{ij} \leq p-1.$$

Substituting from (4.11), we conclude that $tp \in (-p, p)$, so $t \in (-1, 1)$ and hence $t = 0$. It follows that

$$\alpha \left(\sum_{k=1}^m a_{ik} b_{kj} \right) + \beta c_{ij} = d_{ij}.$$

That is, under our assumptions, the congruences (4.8) correctly enforce the desired Hadamard product

$$\alpha(AB) + \beta C = D. \quad (4.12)$$

5. THE MATHS

This section is under construction.

Proposition 5.1. *Let p be a positive integer (not necessarily a prime)*

Proof of Proposition 5.1. □

6. THE CODE

This section is under construction.

Algorithm 6.1 `verify_gen_mat_mul`: verify that $\bar{D} = \bar{\alpha} \cdot \bar{A}\bar{B} + \bar{\beta} \cdot \bar{C}$

Require: Matrices $\bar{A} : \ell \times m, \bar{B} : m \times n, \bar{C}, \bar{D} : \ell \times n$, scalars $\bar{\alpha}, \bar{\beta} \in \mathbb{Z}/p\mathbb{Z}$

Require: All entries are circuit variables representing least residues modulo p

```
1: function VERIFY_GEN_MAT_MUL( $\bar{A}, \bar{B}, \bar{C}, \bar{D}, \bar{\alpha}, \bar{\beta}$ )
2:   for  $i \leftarrow 0$  to  $\ell - 1$  do
3:     for  $j \leftarrow 0$  to  $n - 1$  do
4:        $acc \leftarrow 0$ 
5:       for  $k \leftarrow 0$  to  $m - 1$  do
6:          $acc \leftarrow acc + (\bar{a}_{ik} \cdot \bar{b}_{kj})$ 
7:       end for
8:        $lhs \leftarrow \bar{\alpha} \cdot acc + \bar{\beta} \cdot \bar{c}_{ij}$ 
9:       assert_is_equal( $\bar{d}_{ij}, lhs$ )
10:    end for
11:  end for
12: end function
```

```
1 // Example dimension constants for A, B, C, and D.
2 const L: usize = 3; // Number of rows in A, C, and D
3 const M: usize = 4; // Number of columns in A, rows in B
4 const N: usize = 2; // Number of columns in B, C, and D
5
6 declare_circuit!(Circuit {
7   matrix_a: [[Variable; M]; L], // A: (L x M)
8   matrix_b: [[Variable; N]; M], // B: (M x N)
9   matrix_c: [[Variable; N]; L], // C: (L x N)
10  matrix_d: [[Variable; N]; L], // D = alpha AB + beta C: (L x N)
11  alpha: Value<C::Field>, // scalar alpha (mod p)
12  beta: Value<C::Field>, // scalar beta (mod p)
13});
14
15 impl<C: Config> GenericDefine<C> for Circuit<Variable> {
16   fn define<Builder: RootAPI<C>>(&self, api: &mut Builder) {
17     for i in 0..L {
18       for j in 0..N {
19         let mut acc = api.constant(C::Field::zero());
20
21         for k in 0..M {
22           let prod = api.mul(self.matrix_a[i][k], self.matrix_b[k][j]);
23           acc = api.add(acc, prod);
24         }
25
26         let alpha_acc = api.mul(acc, self.alpha);
27         let beta_c = api.mul(self.matrix_c[i][j], self.beta);
28         let lhs = api.add(alpha_acc, beta_c);
29
30         api.assert_is_equal(self.matrix_d[i][j], lhs);
31       }
32     }
33   }
34 }
```

Listing 1: ECC Rust API: verify $\bar{D} = \bar{\alpha} \cdot \bar{A}\bar{B} + \bar{\beta} \cdot \bar{C}$

7. EXAMPLES

This section is under construction...

REFERENCES

- [1] Inference Labs. *Quantized Matrix Multiplication: Arithmetic Circuit Blueprint*. https://github.com/inference-labs-inc/zkml-blueprints/blob/main/matmul/quantized_matrix_multiplication.pdf. Accessed April 14, 2025.
- [2] Inference Labs. *Range Check and ReLU: Arithmetic Circuit Blueprint*. https://github.com/inference-labs-inc/zkml-blueprints/blob/main/core_ops/range_check_and_relu.pdf. Accessed April 14, 2025.
- [3] Inference Labs. *Range Check, Max, Min, and ReLU: Arithmetic Circuit Blueprint*. https://github.com/inference-labs-inc/zkml-blueprints/blob/main/core_ops/range_check_max_min_relu.pdf. Accessed April 14, 2025.