

Minimal zkVM for Lean Ethereum (draft 0.6.0)

Contents

1	What is the goal of this zkVM?	2
2	VM specification	2
2.1	Field	2
2.2	Memory	3
2.3	Registers	3
2.4	Instruction Set Architecture	3
2.4.1	ADD / MUL	3
2.4.2	DEREF	3
2.4.3	JUMP (Conditional)	4
2.5	Precompiles	4
2.5.1	Poseidon	4
2.5.2	Dot product	4
2.6	ISA programming	4
2.6.1	Functions	4
2.6.2	Loops	4
2.6.3	Range checks	5
2.6.4	Switch statements	5
3	Proving system	6
3.1	Table sizes	6
3.2	Arithmetic Intermediate Representation (AIR)	6
3.3	Execution table	6
3.3.1	Commitment	6
3.3.2	Instruction Encoding	6
3.3.3	AIR transition constraints	7
3.4	Poseidon table	7
3.5	Dot product table	8
3.5.1	Trace layout	8
3.5.2	Columns	8
3.5.3	Memory lookups	9
3.5.4	Bus interaction	9
3.5.5	AIR constraints	9
3.6	Data flow between tables / memory	10
3.6.1	Indexed Lookup into Memory	10
3.6.2	Buses: Data flow between tables	10
3.7	Simple stacking of multilinear polynomials	11
4	Recursive Aggregation	12
4.1	Bytecode evaluation claims	13
4.2	Signer partitioning	13
4.2.1	Trivial partitioning	13
4.2.2	Using a registry of signers	13

1 What is the goal of this zkVM?

Post-quantum signatures are at least an order of magnitude larger than their pre-quantum counterparts, but the migration is necessary to ensure Ethereum’s security. Hash-based signatures offer strong security guarantees and conceptual simplicity, making them a promising candidate: leanXMSS at the consensus layer (where statefulness is not an issue, see [1], [2], and [3]), and leanSPHINCS at the execution layer. A promising choice of hash function is Poseidon2 [4], for its snark-friendliness.

Since post-quantum signatures are much larger, we need to aggregate them. However, they lack the algebraic structure that makes aggregation easy for elliptic-curve schemes like BLS. Instead, we can aggregate them using a snark—itself hash-based, keeping everything post-quantum.

Concretely, we need to:

- **Aggregate** hash-based signatures (leanXMSS / leanSPHINCS)
- **Merge** those aggregate signatures → recursive proof

Both tasks mainly require proving a lot of hashes. A minimal zkVM (inspired by Cairo [5]) is useful as glue to handle all the logic: leanVM.

Aggregation and merging are unified in a single program, that attests that a given set of signers have valid signatures for a given message.

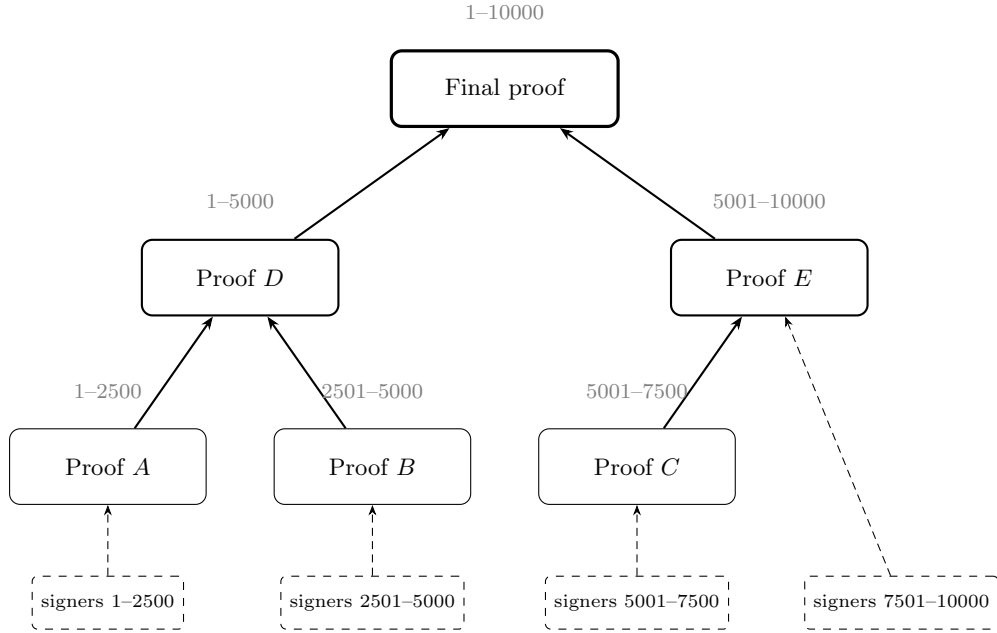


Figure 1: Example of recursive aggregation of 10000 signatures. (Note: overlapping sets of signers are possible.)

2 VM specification

2.1 Field

KoalaBear prime: $p = 2^{31} - 2^{24} + 1$

Advantages:

- small field → less Poseidon rounds
- $x \rightarrow x^3$ is an automorphism of \mathbb{F}_p^* , meaning efficient S-box for Poseidon2 (in BabyBear, it’s degree 7)

- $< 2^{31} \rightarrow$ the sum of 2 field elements can be stored in an u32

The small 2-addicity (24) is not a limiting factor in WHIR, thanks to the use of an interleaved Reed Solomon code (example: with an initial folding of 7, and rate 1/2, we can commit up to 2^{30} field elements).

Extension field: \mathbb{F}_q , with $q = p^5$: the degree-5 extension enables 128 bits of security in WHIR, with the Johnson bound, thanks to the latest result of [6].

2.2 Memory

- Read-Only Memory
- Word = KoalaBear field element
- Memory size: $M = 2^m$ with $16 \leq m \leq 29$ (m depends on the execution and is communicated at the beginning of the proof).
- The first $M' = 2^{m'}$ memory cells hold the "public input", on which both prover and verifier must agree. This enables to pass the arguments that the leanISA program receives as input (in our case: message to sign and XMSS public keys).

2.3 Registers

- pc: program counter
- fp: frame pointer : points to the start of the current stack

Difference with Cairo: no "ap" register (allocation pointer).

2.4 Instruction Set Architecture

Notations:

- α, β and γ represent parameters of the instructions (immediate value operands)
- $\mathbf{m}[i]$ represents the value of the memory at index $i \in \mathbb{F}_p$, with $i < M$ (M : memory size). Any out-of-bound memory access ($i \geq M$) is invalid.
- $\begin{cases} A \\ B \end{cases}$ When using the instruction, either A or B can be used, but not both simultaneously.

2.4.1 ADD / MUL

$a + c = b$ or $a \cdot c = b$ with:

$$a = \begin{cases} \alpha \\ \mathbf{m}[\text{fp} + \alpha] \end{cases} \quad b = \begin{cases} \beta \\ \mathbf{m}[\text{fp} + \beta] \end{cases} \quad c = \begin{cases} \text{fp} \\ \mathbf{m}[\text{fp} + \gamma] \end{cases}$$

2.4.2 Deref

$$\mathbf{m}[\mathbf{m}[\text{fp} + \alpha] + \beta] = \begin{cases} \gamma \\ \mathbf{m}[\text{fp} + \gamma] \\ \text{fp} \end{cases}$$

2.4.3 JUMP (Conditional)

$$\text{condition} = \begin{cases} \alpha \\ \mathbf{m}[\text{fp} + \alpha] \end{cases} \in \{0, 1\} \quad \text{dest} = \begin{cases} \beta \\ \mathbf{m}[\text{fp} + \beta] \end{cases} \quad \text{updated_fp} = \begin{cases} \text{fp} \\ \mathbf{m}[\text{fp} + \gamma] \end{cases}$$

$$\text{next}(\text{pc}) = \begin{cases} \text{dest} & \text{if condition} = 1 \\ \text{pc} + 1 & \text{if condition} = 0 \end{cases} \quad \text{next}(\text{fp}) = \begin{cases} \text{updated_fp} & \text{if condition} = 1 \\ \text{fp} & \text{if condition} = 0 \end{cases}$$

2.5 Precompiles

The ISA supports two precompile instructions, each delegated to a dedicated table in the proving system (Section 3).

2.5.1 Poseidon

TODO

2.5.2 Dot product

Computes $\text{res} = \sum_{i=0}^{n-1} a_i \cdot b_i \in \mathbb{F}_q$, where:

- **BE mode** (`is_be` = 1): $a_i \in \mathbb{F}_p$, $b_i \in \mathbb{F}_q$ (base \times extension). The a_i are read from consecutive memory addresses starting at α , and the b_i from consecutive blocks of 5 starting at β .
- **EE mode** (`is_be` = 0): $a_i \in \mathbb{F}_q$, $b_i \in \mathbb{F}_q$ (extension \times extension). Both slices are read from consecutive blocks of 5, starting at α and β respectively.

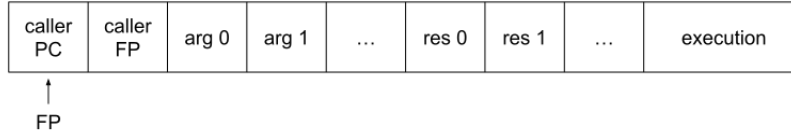
Parameters: memory addresses α (first operand), β (second operand), γ (result), a size $n > 0$, and the mode flag `is_be` $\in \{0, 1\}$. The result is written as 5 consecutive field elements at $\mathbf{m}[\gamma.. \gamma + 5]$.

2.6 ISA programming

2.6.1 Functions

1. Each function has a deterministic memory footprint: the length of the continuous frame in memory that is allocated for each of its calls.
2. At runtime, each time we call our function, we receive via a memory cell a hint pointing to a free range of memory. We then store the current values of pc / fp at the start of this newly allocated frame, alongside the function's arguments, we can then jump, to enter the function bytecode, and modify fp with the hinted value. The intuition here is that the verifier does not care where the new memory frame will be placed (we use a read-only memory, so we cannot overwrite previous frames). In practice, the prover that runs the program would need to keep the value of the allocation pointer "ap", in order to adequately allocate new memory frames, but there is no need to keep track of it from the verifier's perspective.

Figure 2: Memory layout of a function call



2.6.2 Loops

We suggest to unroll loops when the number of iterations is low, and known at compile time. The remaining loops are transformed into recursive functions (by the leanISA compiler).

2.6.3 Range checks

It's possible to check that a given memory cell is smaller than some value t (for $t \leq 2^{16}$) in 3 cycles.

We denote by $\mathbf{m}[\text{fp} + x]$ the memory cell for which we want to ensure $\mathbf{m}[\text{fp} + x] < t$. We also denote by $\mathbf{m}[\text{fp} + i]$, $\mathbf{m}[\text{fp} + j]$ and $\mathbf{m}[\text{fp} + k]$ 3 auxiliary memory cells (that have not been used yet).

1. $\mathbf{m}[\mathbf{m}[\text{fp} + x]] = \mathbf{m}[\text{fp} + i]$ (using DEREf, this ensures $\mathbf{m}[\text{fp} + x] < M$, the memory size)
2. $\mathbf{m}[\mathbf{m}[\text{fp} + x]] + \mathbf{m}[\text{fp} + j] = (t-1)$ (using ADD)
3. $\mathbf{m}[\mathbf{m}[\text{fp} + j]] = \mathbf{m}[\text{fp} + k]$ (using DEREf, this ensures $t - 1 - \mathbf{m}[\text{fp} + x] < M$)

Given $t \leq 2^{16} \leq M$, $\mathbf{m}[\text{fp} + x] < M$, $t - 1 - \mathbf{m}[\text{fp} + x] < M$, and $M \leq 2^{29} < p/2$, we have: $\mathbf{m}[\text{fp} + x] < t$.

Note: From the point of view of the prover running the program, some hints are necessary (filling the values of $\mathbf{m}[\text{fp} + i]$ and $\mathbf{m}[\text{fp} + k]$ must be done at end of execution).

This idea was pointed out by D. Khovratovich, and is an unplanned use of the DEREf instruction.

Example 2.1. Let's say we want to write a function with 2 arguments $x = \mathbf{m}[\text{fp} + 2]$ and $y = \mathbf{m}[\text{fp} + 3]$ ($\mathbf{m}[\text{fp} + 0]$ and $\mathbf{m}[\text{fp} + 1]$ are used, by convention, to store the caller's pc and fp, to return to the previous context at the end of the function), which perform the following:

1. `assert(x < 10)`
2. `z := x*y + 100`
3. `assert(z < 1000)`

Which can be compiled to:

1. $\mathbf{m}[\text{fp} + 4] = \mathbf{m}[\text{fp} + 2]$ // check that x is "small"
2. $\mathbf{m}[\text{fp} + 2] + \mathbf{m}[\text{fp} + 5] = 9$ // compute $9 - x$
3. $\mathbf{m}[\text{fp} + 6] = \mathbf{m}[\text{fp} + 5]$ // check that $9 - x$ is "small"
4. $\mathbf{m}[\text{fp} + 7] = \mathbf{m}[\text{fp} + 2] * \mathbf{m}[\text{fp} + 3]$ // compute $x.y$
5. $\mathbf{m}[\text{fp} + 8] = \mathbf{m}[\text{fp} + 7] + 100$ // compute $z = x.y + 100$
6. $\mathbf{m}[\text{fp} + 9] = \mathbf{m}[\text{fp} + 8]$ // check that z is "small"
7. $\mathbf{m}[\text{fp} + 8] + \mathbf{m}[\text{fp} + 10] = 999$ // compute $999 - z$
8. $\mathbf{m}[\text{fp} + 11] = \mathbf{m}[\text{fp} + 10]$ // check that $999 - z$ is "small"
9. JUMP with `next(pc) = $\mathbf{m}[\text{fp} + 0]$` , `next(fp) = $\mathbf{m}[\text{fp} + 1]$` , `condition = 1` // return

2.6.4 Switch statements

Suppose we want a different logic depending on the value x of a given memory cell, where x is known to be $< k$ (if the value x comes from a "hint", don't forget to range-check it).

Each of the k different value leads to a different branch at runtime, represented by a block of code. We want to jump to the correct block of code depending on x . One efficient implementation consists in placing our blocks of code at regular intervals, and to jump to a $a + b.x$, where a is the offset of the first block of code (in case $x = 0$), and b is the distance between two consecutive blocks.

Example: During XMSS verification, for each of the v chains, we need to hash a pre-image, a number of times depending on the encoding, but known to be $< w$. Here $k = w$, and the i -th block of code we could jump to will execute i times the hash function (unrolled loop).

3 Proving system

3.1 Table sizes

Each table has a maximum number of rows (as a power of two):

Table	Max rows
Execution	2^{25}
Dot product	2^{21}
Poseidon	2^{21}

Tables are padded to the next power of two (with a minimum of 2^8 rows). These bounds are critical for soundness: they ensure that logup multiplicities cannot overflow modulo p .

3.2 Arithmetic Intermediate Representation (AIR)

If a table has n columns, a transition constraint is a multivariate polynomial $C(x_1, \dots, x_n, y_1, \dots, y_n)$ in $2n$ variables, where (x_1, \dots, x_n) represents the current row and (y_1, \dots, y_n) represents the next row. Each table defines a set of such constraints, all of which must evaluate to zero on every pair of consecutive rows. For a table with H rows (indexed $0, \dots, H-1$), the constraints are evaluated on $(H-1) + 1 = H$ pairs:

- $H-1$ consecutive pairs: $(\text{row}_i, \text{row}_{i+1})$ for $i = 0, \dots, H-2$
- 1 **wrapping pair**: $(\text{row}_{H-1}, \text{row}_{H-1})$ — the last row is paired with itself

The wrapping pair means that on the last row, $\text{next}(x) = x$ for every column x .

TODO: describe how to prove AIR via sumcheck.

3.3 Execution table

3.3.1 Commitment

At each cycle, we commit to 20 (base) field elements:

- pc (program counter)
- fp (frame pointer)
- $\text{addr}_A, \text{addr}_B, \text{addr}_C$
- $\text{value}_A = \mathbf{m}[\text{addr}_A], \text{value}_B = \mathbf{m}[\text{addr}_B], \text{value}_C = \mathbf{m}[\text{addr}_C]$
- 12 field elements describing the instruction being executed (see below)

3.3.2 Instruction Encoding

Each instruction is described by 12 field elements:

- 3 operands ($\in \mathbb{F}_p$): $\text{operand}_A, \text{operand}_B, \text{operand}_C$
- 3 associated flags ($\in \{0, 1\}$): $\text{flag}_A, \text{flag}_B, \text{flag}_C$
- 5 opcode flags ($\in \{0, 1\}$): ADD, MUL, Deref, JUMP, PRECOMPILE_INDEX
- 1 multi-purpose operand: AUX

3.3.3 AIR transition constraints

We use transition constraints of degree 5, but it's always possible to make them quadratic with additional columns in the execution table.

We define the following quantities:

- $\nu_A = \text{flag}_A \cdot \text{operand}_A + (1 - \text{flag}_A) \cdot \text{value}_A$
- $\nu_B = \text{flag}_B \cdot \text{operand}_B + (1 - \text{flag}_B) \cdot \text{value}_B$
- $\nu_C = \text{flag}_C \cdot \text{fp} + (1 - \text{flag}_C) \cdot \text{value}_C$

With the associated constraints: $\forall X \in \{A, B, C\} : (1 - \text{flag}_X) \cdot (\text{address}_X - (\text{fp} + \text{operand}_X)) = 0$

For addition and multiplication:

- $\text{ADD} \cdot (\nu_B - (\nu_A + \nu_C)) = 0$
 - $\text{MUL} \cdot (\nu_B - \nu_A \cdot \nu_C) = 0$
-

When $\text{DEREF} = 1$, set $\text{flag}_A = 0$, $\text{flag}_C = 1$ and:

$$\mathbf{m}[\mathbf{m}[\text{fp} + \alpha] + \beta] = \begin{cases} \gamma & \rightarrow \text{AUX} = 1, \text{flag}_B = 1 \\ \mathbf{m}[\text{fp} + \gamma] & \rightarrow \text{AUX} = 1, \text{flag}_B = 0 \\ \text{fp} & \rightarrow \text{AUX} = 0 (\text{flag}_B = 1) \end{cases}$$

- $\text{DEREF} \cdot (\text{addr}_C - (\text{value}_A + \text{operand}_C)) = 0$
 - $\text{DEREF} \cdot \text{AUX} \cdot (\text{value}_C - \nu_B) = 0$
 - $\text{DEREF} \cdot (1 - \text{AUX}) \cdot (\text{value}_C - \text{fp}) = 0$
-

For jump and conditional state transitions, let $J = \text{JUMP} \cdot \nu_A$:

- $\text{JUMP} \cdot \nu_A \cdot (1 - \nu_A) = 0$
- $J \cdot (\text{next}(\text{pc}) - \nu_B) = 0$
- $J \cdot (\text{next}(\text{fp}) - \nu_C) = 0$
- $(1 - J) \cdot (\text{next}(\text{pc}) - (\text{pc} + 1)) = 0$
- $(1 - J) \cdot (\text{next}(\text{fp}) - \text{fp}) = 0$

Note: the constraint $\text{JUMP} \cdot \nu_A \cdot (1 - \nu_A) = 0$ could be removed, as long as it's correctly enforced in the bytecode.

TODO: Verify and (formally) prove soundness.

3.4 Poseidon table

TODO

3.5 Dot product table

The dot product precompile computes $\text{res} = \sum_{i=0}^{n-1} a_i \cdot b_i \in \mathbb{F}_q$, where:

- **BE mode** (`is_be = 1`): $a_i \in \mathbb{F}_p$, $b_i \in \mathbb{F}_q$ (base \times extension)
- **EE mode** (`is_be = 0`): $a_i \in \mathbb{F}_q$, $b_i \in \mathbb{F}_q$ (extension \times extension)

3.5.1 Trace layout

Each dot product of length n occupies n consecutive rows. The first row (`start = 1`) carries the full accumulated result; subsequent rows (`start = 0`) count down via `len`. Multiple dot products are stacked vertically.

is_be	flag	Base field columns (\mathbb{F}_p)						Extension field columns (\mathbb{F}_q)				aux
		start	len	idx _A	idx _B	idx _R	v_{A,\mathbb{F}_p}	v_{A,\mathbb{F}_q}	v_B	res	comp	
<i>BE dot product</i> ($n = 3, \alpha = 90, \beta = 200, \gamma = 50$):								$\text{res} = \sum_{i=0}^2 \mathbf{m}[90+i] \cdot \mathbf{m}[200+5i \dots 205+5i]$				
1	1	1	3	90	200	50	$\mathbf{m}[90]$	$\mathbf{m}[90..95]$	$\mathbf{m}[200..205]$	$\mathbf{m}[50..55]$	c_2	7
1	0	0	2	91	205	50	$\mathbf{m}[91]$	$\mathbf{m}[91..96]$	$\mathbf{m}[205..210]$	$\mathbf{m}[50..55]$	c_1	5
1	0	0	1	92	210	50	$\mathbf{m}[92]$	$\mathbf{m}[92..97]$	$\mathbf{m}[210..215]$	$\mathbf{m}[50..55]$	c_0	3
<i>EE dot product</i> ($n = 2, \alpha = 300, \beta = 400, \gamma = 500$):								$\text{res} = \sum_{i=0}^1 \mathbf{m}[300+5i \dots 305+5i] \cdot \mathbf{m}[400+5i \dots 405+5i]$				
0	1	1	2	300	400	500	$\mathbf{m}[300]$	$\mathbf{m}[300..305]$	$\mathbf{m}[400..405]$	$\mathbf{m}[500..505]$	d_1	4
0	0	0	1	305	405	500	$\mathbf{m}[305]$	$\mathbf{m}[305..310]$	$\mathbf{m}[405..410]$	$\mathbf{m}[500..505]$	d_0	2

Recurrence (over \mathbb{F}_q):

- BE: $c_0 = \mathbf{m}[92] \cdot \mathbf{m}[210..215]$, $c_1 = \mathbf{m}[91] \cdot \mathbf{m}[205..210] + c_0$, $c_2 = \mathbf{m}[90] \cdot \mathbf{m}[200..205] + c_1 = \text{res}$
- EE: $d_0 = \mathbf{m}[305..310] \cdot \mathbf{m}[405..410]$, $d_1 = \mathbf{m}[300..305] \cdot \mathbf{m}[400..405] + d_0 = \text{res}$

Notes: In BE mode, `idxA` increments by 1 (base field elements); in EE mode, by 5 (extension field elements). `idxB` always increments by 5. `aux = is_be + 2 · len` is a non-committed column (used only for bus fingerprinting). Both v_{A,\mathbb{F}_p} and v_{A,\mathbb{F}_q} are read from memory on every row, but only the relevant one (depending on `is_be`) is used in the computation.

3.5.2 Columns

8 base field columns + 4 extension field columns = 28 \mathbb{F}_p -columns.

Base field columns (\mathbb{F}_p):

- `is_be` $\in \{0, 1\}$: mode flag (1 = base \times extension, 0 = extension \times extension)
- `flag` $\in \{0, 1\}$: bus selector (1 only on the first row of a real (non-padding) dot product)
- `start` $\in \{0, 1\}$: 1 on the first row of each group (and on padding rows)
- `len`: countdown from n to 1
- `idxA`, `idxB`, `idxR`: memory addresses for A , B , and the result
- v_{A,\mathbb{F}_p} : base field value $\mathbf{m}[\text{idx}_A]$ (used in BE mode)

Extension field columns (\mathbb{F}_q):

- v_{A,\mathbb{F}_q} : extension field value $\mathbf{m}[\text{idx}_A..\text{idx}_A + 5]$
- v_B : extension field value $\mathbf{m}[\text{idx}_B..\text{idx}_B + 5]$
- `res`: result value $\mathbf{m}[\text{idx}_R..\text{idx}_R + 5]$
- `comp`: running computation (accumulated dot product)

Non-committed column (used for bus fingerprinting):

- `aux = is_be + 2 · len`

3.5.3 Memory lookups

Every row performs:

- 1 base field lookup: $v_{A,\mathbb{F}_p} = \mathbf{m}[\text{idx}_A]$
- 3 extension field lookups: $v_{A,\mathbb{F}_q} = \mathbf{m}[\text{idx}_A.. + 5]$, $v_B = \mathbf{m}[\text{idx}_B.. + 5]$, $\text{res} = \mathbf{m}[\text{idx}_R.. + 5]$

3.5.4 Bus interaction

On rows where $\text{flag} = 1$, the table PULLS $(\text{idx}_A, \text{idx}_B, \text{idx}_R, \text{aux})$ from the corresponding bus. The execution table PUSHes a matching tuple for each dot product instruction.

The aux encoding ensures both the mode and the length are bound to the bus data. Since is_be is constrained to be boolean (constraint 3), if an adversary flips is_be while keeping aux constant, the adversary would need a length len' such that $\text{aux} = (1 - \text{is_be}) + 2 \cdot \text{len}'$, i.e. $2 \cdot \text{len}' = 2 \cdot \text{len} \pm 1$ over \mathbb{F}_p . By Lemma 1, both len and len' are $\leq 2^{21} < p/2$, so the equation cannot overflow modulo p and must hold over \mathbb{Z} . Over \mathbb{Z} , this is impossible by parity.

3.5.5 AIR constraints

We write $\text{next}(x)$ for the value of column x on the **next row** (down).

Value selection: Define

$$v_A = \text{is_be} \cdot v_{A,\mathbb{F}_p} + (1 - \text{is_be}) \cdot v_{A,\mathbb{F}_q} \in \mathbb{F}_q$$

Constraints:

1. $\text{start} \cdot (1 - \text{start}) = 0$ (start is boolean)
2. $\text{flag} \cdot (1 - \text{flag}) = 0$ (flag is boolean)
3. $\text{is_be} \cdot (1 - \text{is_be}) = 0$ (is_be is boolean)
4. $\text{flag} \cdot (1 - \text{start}) = 0$ ($\text{flag} \Rightarrow \text{start}$)
5. $(\text{next}(\text{is_be}) - \text{is_be})^2 \cdot (1 - \text{next}(\text{start})) = 0$ (mode consistent within group)
6. $\text{comp} - v_A \cdot v_B - \text{next}(\text{comp}) \cdot (1 - \text{next}(\text{start})) = 0$ (recurrence, over \mathbb{F}_q)
7. $(1 - \text{next}(\text{start})) \cdot (\text{len} - \text{next}(\text{len}) - 1) = 0$ (length countdown)
8. $\text{next}(\text{start}) \cdot (\text{len} - 1) = 0$ (must reach 1 before new group)
9. $(1 - \text{next}(\text{start})) \cdot (\text{idx}_A - \text{next}(\text{idx}_A) + \delta_A) = 0$ (with $\delta_A = 5(1 - \text{is_be}) + \text{is_be}$)
10. $(1 - \text{next}(\text{start})) \cdot (\text{idx}_B - \text{next}(\text{idx}_B) + 5) = 0$
11. $(\text{comp} - \text{res}) \cdot \text{start} = 0$ (result check, over \mathbb{F}_q)

Constraint 6: when $\text{next}(\text{start}) = 1$ (next row starts a new group), the tail term vanishes: $\text{comp} = v_A \cdot v_B$. Otherwise the partial sum is accumulated. Constraint 11: on the first row of each group, the accumulated computation must equal the result read from memory.

TODO: Verify and (formally) prove soundness.

Lemma 1. On every row r of the dot product table, $\text{len}_r \leq H - r$ (as an integer in $\{0, \dots, p-1\}$), where $H \leq 2^{21}$ is the table height. In particular, $\text{len}_r \leq H$.

Proof. By backward induction on r , from $r = H - 1$ down to $r = 0$.

Base case ($r = H - 1$): On the wrapping pair (row $H - 1$ paired with itself, see Section 3.2), $\text{next}(x) = x$ for every column x . Constraint 7 becomes $(1 - \text{start}) \cdot (\text{len} - \text{len} - 1) = -(1 - \text{start}) = 0$, forcing $\text{start}_{H-1} = 1$. Then constraint 8 becomes $1 \cdot (\text{len} - 1) = 0$, giving $\text{len}_{H-1} = 1 = H - (H - 1)$.

Inductive step ($r < H - 1$): Assume $\text{len}_s \leq H - s$ for all $s > r$. The pair (row r , row $r+1$) gives two cases depending on $\text{start}_{r+1} \in \{0, 1\}$ (constraint 1):

- If $\text{start}_{r+1} = 1$: constraint 8 gives $1 \cdot (\text{len}_r - 1) = 0$, so $\text{len}_r = 1 \leq H - r$.
- If $\text{start}_{r+1} = 0$: constraint 7 gives $(1 - 0) \cdot (\text{len}_r - \text{len}_{r+1} - 1) = 0$, so $\text{len}_r = \text{len}_{r+1} + 1$ in \mathbb{F}_p . By induction, $\text{len}_{r+1} \leq H - (r + 1)$, so $\text{len}_{r+1} + 1 \leq H - r$. Since $H - r \leq H \leq 2^{21} < p$, the addition does not overflow modulo p , so $\text{len}_r = \text{len}_{r+1} + 1 \leq H - r$ as integers.

□

3.6 Data flow between tables / memory

Lemma 2. Let a_0, a_1, \dots, a_{n-1} be pairwise distinct poles in \mathbb{F}_q , and let m_0, m_1, \dots, m_{n-1} be an associated list of multiplicities in $\{0, 1, \dots, p - 1\}$. Consider the rational function:

$$P(X) = \sum_{i=0}^{n-1} \frac{m_i}{X - a_i}$$

Except with probability n/q , if $P(\alpha) = 0$ for a random $\alpha \in \mathbb{F}_q$, then all multiplicities $m_i = 0$.

3.6.1 Indexed Lookup into Memory

We use `logup` [7], in its indexed form, to allow tables to perform lookups into the read-only memory.

Let \mathcal{T} denote the set of all tables in the system. For each table $T \in \mathcal{T}$ with H_T rows, let n_T denote the number of memory lookups. Each lookup $i < n_T$ consists of an **index column** $\text{col}_{\text{index}, T, i}$ and a **value column** $\text{col}_{\text{val}, T, i}$.

The rule to enforce is the following:

$$\forall T \in \mathcal{T}, \forall i < n_T, \forall j < H_T : \quad \text{col}_{\text{val}, T, i}(j) = \mathbf{m}[\text{col}_{\text{index}, T, i}(j)]$$

Implicitly, we must also have $\text{col}_{\text{index}, T, i}(j) < M$ (the memory size).

The prover initially commits to a multilinear polynomial acc , having the same size as the memory, such that (in the honest case) for every $k < M$:

$$\text{acc}[k] = \sum_{T \in \mathcal{T}} \sum_{i < n_T} |\{j < H_T : \text{col}_{\text{index}, T, i}(j) = k\}|$$

i.e., $\text{acc}[k]$ represents the total number of times address k is accessed by the lookups across all tables.

The verifier sends a random challenge $\alpha \in \mathbb{F}_q$ (TODO quantify soundness error). Let $N = \sum_{T \in \mathcal{T}} n_T \cdot H_T$ be the total number of memory lookups. Assuming $N < p$ (to avoid overflow), the indexed lookup into memory is valid, except with probability $(N + M)/q$, if for a randomly sampled $X \in \mathbb{F}_q$:

$$\sum_{T \in \mathcal{T}} \sum_{i < n_T} \sum_{j < H_T} \frac{1}{X - (\text{col}_{\text{index}, T, i}(j) + \alpha \cdot \text{col}_{\text{val}, T, i}(j))} = \sum_{k < M} \frac{\text{acc}(k)}{X - (k + \alpha \cdot \mathbf{m}[k])}$$

This can be computed via GKR, as introduced in [8].

3.6.2 Buses: Data flow between tables

See OpenVM [9], definition 2.2.2, for more details.

To each table is associated a list of buses. Declaring a bus for a given table requires specifying:

1. **Direction**: either PULL or PUSH
2. **Bus index**: can be either a constant or given by a column
3. **Data columns**: a list (of arbitrary size) of columns defining which data is passed to the bus

4. **Selector column:** enforced via AIR constraint to only contain Boolean values ($\in \{0, 1\}$), indicating at which row the bus is active (1) or silenced (0)

Balance rule: At the end of execution, every bus (identified by its index) must be *balanced* across all tables: for each data tuple, the number of pushes must equal the number of pulls.

Proving system (logup):

Let \mathcal{T} denote the set of all tables in the system. For each table $T \in \mathcal{T}$ with H_T rows, let \mathcal{B}_T denote the set of buses associated with T . For each bus $b \in \mathcal{B}_T$, we have:

- $\text{dir}_{T,b} \in \{+1, -1\}$: the direction (+1 for PUSH, -1 for PULL)
- $\text{sel}_{T,b} : [0, H_T) \rightarrow \{0, 1\}$: the selector column
- $\text{bus_idx}_{T,b} : [0, H_T) \rightarrow \mathbb{F}_p$: the bus index (assumed to be a column here, but could also be a constant)
- $\text{data}_{T,b} = (d_0, d_1, \dots, d_{k-1})$: the k data columns

The verifier sends a random challenge $\alpha \in \mathbb{F}_q$ (TODO quantify soundness error). For a bus with index β and data tuple $(c_0, c_1, \dots, c_{k-1})$, define the encoding:

$$\text{encode}(\beta, c_0, \dots, c_{k-1}) = \beta + \sum_{i=0}^{k-1} \alpha^{i+1} \cdot c_i$$

Except with probability N/q , where $N = \sum_{T \in \mathcal{T}} |\mathcal{B}_T| \cdot H_T$ is the total number of bus interactions, the bus balance constraint is satisfied if for a randomly sampled $X \in \mathbb{F}_q$:

$$\sum_{T \in \mathcal{T}} \sum_{b \in \mathcal{B}_T} \sum_{j=0}^{H_T-1} \frac{\text{dir}_{T,b} \cdot \text{sel}_{T,b}(j)}{X - \text{encode}(\text{bus_idx}_{T,b}(j), \text{data}_{T,b}(j))} = 0$$

Overflow constraint: For soundness, we must ensure that multiplicities do not overflow modulo p . This requires that for each bus index β , the total number of interactions (pushes and pulls combined) is strictly less than p :

$$\forall \beta : \sum_{\substack{T \in \mathcal{T}, b \in \mathcal{B}_T \\ \text{bus_idx}_{T,b} = \beta}} H_T < p$$

As with the indexed lookup into memory, this sum can be computed efficiently via GKR [8]. In practice, all logup sums (bus balance and memory lookups) are batched into a single GKR instance, using random linear combination.

TODO: Unified vision in which memory is just another table, and memory accesses are just bus interactions?

A detailed soundness analysis can be found in [Soundness of Interactions via LogUp](#).

3.7 Simple stacking of multilinear polynomials

Note 1: It's always possible to reduce n claims about a multilinear polynomial to a single one, using sumcheck. But this trick is not necessary with WHIR, which natively supports an arbitrary number of claims about the committed polynomial.

Note 2: Crucially, the proving cost to add an equality constraint of the form $P((x_1, \dots, x_n)) = y$ to a committed polynomial P via WHIR is $O(2^k)$ where $k = |\{i, x_i \notin \{0, 1\}\}|$ is the number of "non-boolean variables". As a result, adding "sparse" $((x_i)$ containing boolean values) equality constraints is essentially optimal.

In order to commit to multiple univariate polynomials with FRI, each polynomial must be FFT-ed + Merkle-committed. Even if it's possible to have some batching at the Merkle tree level (see 'MMCS' in [Plonky3](#)), the proof size for multiple, complex AIR tables quickly reach the megabyte scale.

With a multilinear PCS (such as WHIR), we can "concatenate" multiple multilinear polynomials into a single one, and commit to it once (offering significant proof size savings).

More details: Given n multilinear polynomials P_1, \dots, P_n with ν_1, \dots, ν_n variables respectively, we order them from the **largest to the smallest** and concatenate their respective evaluation (on the boolean hypercube):

$$P = [P_1(\{0, 1\}^{\nu_1}) \| P_2(\{0, 1\}^{\nu_2}) \| \dots \| P_n(\{0, 1\}^{\nu_n})]$$

After padding with zeros to the next power of two, we can interpret the result as the evaluations (on the boolean hypercube) of a multilinear polynomial, with $\nu = \lceil \log_2(\sum_i 2^{\nu_i}) \rceil$ variables, and commit to it.

To reduce an evaluation claim on an "inner" (smaller) polynomial P_i to a claim on the "outer" (larger) polynomial P , we use **boolean selectors**.

Example: Consider 3 multilinear polynomials P_1, P_2, P_3 with 4, 3, and 2 variables respectively. The concatenated polynomial P , that we commit, has $5 = \lceil \log_2(2^4 + 2^3 + 2^2) \rceil$ variables.

- $P_1(x_1, x_2, x_3, x_4) = P(0, x_1, x_2, x_3, x_4)$
- $P_2(x_1, x_2, x_3) = P(1, 0, x_1, x_2, x_3)$
- $P_3(x_1, x_2) = P(1, 1, 0, x_1, x_2)$

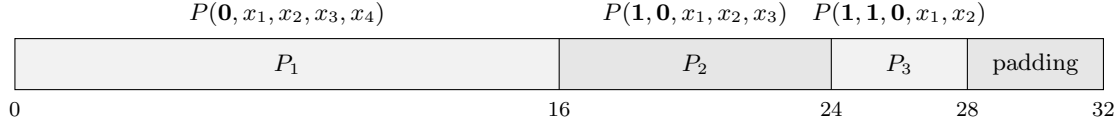


Figure 3: Simple stacking of P_1, P_2, P_3 into a single polynomial P

Advantage of this approach: simplicity.

Drawback: padding overhead, i.e. it does not take advantage of the potential repetitions at the end of each inner polynomial's evaluations.

There are alternative ways to handle the stacking of multiple multilinear polynomials:

- **Jagged PCS** [10]: No padding overhead, at the cost of an additional sumcheck.
- **Per-polynomial chunking**: decompose the non-repeated part of each inner polynomial into a small (think 3 or 4) number of power-of-two sized chunks, and concatenate all these chunks into a single large polynomial.

Note 3: A meticulous implementation of WHIR could also take advantage of any potential repetitions in the committed polynomial, both at the FFT and the sumcheck level, leaving the Merkleisation as the main overhead.

4 Recursive Aggregation

The recursive aggregation program (see Figure 1) takes n_{signs} signers and partitions them into $n_{\text{rec}} + 1$ sources: one **direct source** whose signatures are verified explicitly, and n_{rec} **recursive sources**, each accompanied by a sub-proof that is recursively verified.

For each recursive source, the program runs the leanVM verifier on the provided sub-proof, confirming that the sub-proof itself is a valid recursive aggregation proof covering the source's signers.

4.1 Bytecode evaluation claims

A key subtlety arises from recursion. Every leanVM proof is tied to a specific *bytecode* — the program that was executed, represented as a multilinear polynomial. In our case, there is a single program (that both contain signature verification and recursion algorithms). The leanVM verification algorithm requires evaluating the bytecode polynomial at a random point.

When the recursive aggregation program verifies a sub-proof, it runs the leanVM verifier *inside* the program. This verification produces a bytecode evaluation claim. For efficiency, rather than checking it inside the program (with a PCS), the claim is forwarded to the public input, to be checked externally.

However, each sub-proof is potentially itself a recursive aggregation, and may have already forwarded its own bytecode claim in the same way. So for each of the n_{rec} sub-proofs, **two** bytecode claims arise:

1. **Inner-public-memory claim:** The bytecode claim that the sub-proof forwarded from its own recursive verifications (read from the sub-proof’s public input).
2. **Inner-proof claim:** The new bytecode claim produced by verifying the sub-proof itself.

After verifying all n_{rec} sub-proofs, a fresh Fiat–Shamir transcript is initialized and fed with all $2n_{\text{rec}}$ evaluation points and claimed values. This transcript produces a random linear combination challenge, and the $2n_{\text{rec}}$ claims are then batched into a single one via sumcheck, yielding a reduced claim. This reduced claim is written to the public input for the next level of recursion.

When $n_{\text{rec}} = 0$ (no recursion), the bytecode claim in the public input is set to zero.

4.2 Signer partitioning

4.2.1 Trivial partitioning

In the most straightforward approach, the aggregation program takes as input a list of public keys, and a message. In a loop, each public key is either directly verified (direct source) or assigned to one of the n_{rec} recursive sources. Then, the program verifies each recursive source accordingly (with the same message).

- Pros: simplicity
- Cons: performance bottleneck in the final steps of aggregation when there are a lot of signers (think more than 2^{14})

If performance allows, this approach should be preferred.

4.2.2 Using a registry of signers

When the set of signers is large, but is known in advance (specifically: Ethereum’s consensus layer), a potentially more efficient approach is possible, yet more complex.

All public keys are stored in a Merkle tree (the *public-key registry*). Signers are identified by their *index* in this registry (rather than using a bitfield).

For the direct source, the program checks each signature and verifies a Merkle inclusion proof that their public key belongs to the registry.

Each aggregated signature attests to a set of n_{sig} signer indices from the public-key registry. The $n_{\text{rec}} + 1$ sources (one direct source and n_{rec} recursive sources) may *overlap*: a signer can appear in multiple sources. The program must verify that the claimed global set included in the union of all sources.

The (untrusted) prover hints three kinds of index lists:

- **Global indices G :** the deduplicated list of all signer indices attested by this proof.
- **Duplicate indices D :** each index appearing in $k > 1$ sources contributes $k - 1$ copies to D .
- **Source indices S_s** for each source $s \in \{0, \dots, n_{\text{rec}}\}$.

A constraint asserts $|G| + |D| = \sum_{s=0}^{n_{\text{rec}}} |S_s|$.
The multiset identity to enforce is:

$$G \cup D = \bigcup_{s=0}^{n_{\text{rec}}} S_s$$

Index hashing. All these $2 + n_{\text{rec}}$ lists are hashed: $H(G)$, $H(D)$, and $(H(S_s))_{s=0}^{n_{\text{rec}}}$. $H(G)$ is written to the public input (outer proof), and each $H(S_s)$ is written to the public input of the corresponding inner proof.

Challenge derivation. By hashing together $H(G)$, $H(D)$, and $(H(S_s))_{s=0}^{n_{\text{rec}}}$, a Fiat–Shamir challenge $\alpha \in \mathbb{F}_q$ is derived.

Consistency checks. For any index list $I = (i_0, \dots, i_{k-1})$ with $i_j \in \mathbb{F}_p$, define the product polynomial $P_I(\alpha) = \prod_{j=0}^{k-1} (\alpha - i_j) \in \mathbb{F}_q$. The program evaluates both sides and asserts:

$$P_G(\alpha) \cdot P_D(\alpha) = \prod_{s=0}^{n_{\text{rec}}} P_{S_s}(\alpha)$$

Soundness of this argument is $(|G| + |D|)/q$, by the Schwartz–Zippel lemma.

The current implementation uses this signature registry, TODO can we use the trivial partitioning instead (without too much performance overhead)?

References

- [1] J. Drake, D. Khovratovich, M. Kudinov, and B. Wagner, “Hash-based multi-signatures for post-quantum ethereum,” Cryptology ePrint Archive, Paper 2025/055, 2025. [Online]. Available: <https://eprint.iacr.org/2025/055>
- [2] D. Khovratovich, M. Kudinov, and B. Wagner, “At the top of the hypercube – better size-time tradeoffs for hash-based signatures,” Cryptology ePrint Archive, Paper 2025/889, 2025. [Online]. Available: <https://eprint.iacr.org/2025/889>
- [3] J. Drake, D. Khovratovich, M. Kudinov, and B. Wagner, “Technical note: LeanSig for post-quantum ethereum,” Cryptology ePrint Archive, Paper 2025/1332, 2025. [Online]. Available: <https://eprint.iacr.org/2025/1332>
- [4] L. Grassi, D. Khovratovich, and M. Schofnegger, “Poseidon2: A faster version of the poseidon hash function,” Cryptology ePrint Archive, Paper 2023/323, 2023. [Online]. Available: <https://eprint.iacr.org/2023/323>
- [5] L. Goldberg, S. Papini, and M. Riabzev, “Cairo – a turing-complete STARK-friendly CPU architecture,” Cryptology ePrint Archive, Paper 2021/1063, 2021. [Online]. Available: <https://eprint.iacr.org/2021/1063>
- [6] E. Ben-Sasson, D. Carmon, U. Haböck, S. Kopparty, and S. Saraf, “On proximity gaps for reed-solomon codes,” Cryptology ePrint Archive, Paper 2025/2055, 2025. [Online]. Available: <https://eprint.iacr.org/2025/2055>
- [7] U. Haböck, “Multivariate lookups based on logarithmic derivatives,” Cryptology ePrint Archive, Paper 2022/1530, 2022. [Online]. Available: <https://eprint.iacr.org/2022/1530>
- [8] S. Papini and U. Haböck, “Improving logarithmic derivative lookups using GKR,” Cryptology ePrint Archive, Paper 2023/1284, 2023. [Online]. Available: <https://eprint.iacr.org/2023/1284>
- [9] O. CONTRIBUTORS, “Openvm whitepaper,” 2025. [Online]. Available: <https://openvm.dev/whitepaper.pdf>
- [10] T. Hemo, K. Jue, E. Rabinovich, G. Roh, and R. D. Rothblum, “Jagged polynomial commitments (or: How to stack multilinear),” Cryptology ePrint Archive, Paper 2025/917, 2025. [Online]. Available: <https://eprint.iacr.org/2025/917>