

Folding-based zkLLM

Wilbert Wu

Abstract

This paper introduces a new approach to construct zero-knowledge large language models (zkLLM) based on the Folding technique. We first review the concept of Incrementally Verifiable Computation (IVC) and compare the IVC constructions based on SNARK and Folding. Then we discuss the necessity of Non-uniform IVC (NIVC) and present several Folding schemes that support more expressive circuits, such as SuperNova, Sangria, Origami, HyperNova, and Protostar. Based on these techniques, we propose a zkLLM design that uses a RAM machine architecture with a set of opcodes. We define corresponding constraint circuits for each opcode and describe the workflows of the prover and verifier. Finally, we provide examples of opcodes to demonstrate the circuit construction methods. Our zkLLM design achieves high efficiency and expressiveness, showing great potential for practical applications.

1 IVC

IVC (Incrementally Verifiable Computation) is a new primitive in Folding ZK. Based on Folding, IVC can be constructed more efficiently, where the prover proves the correctness of incremental computation $y = F^{(n)}(x)$.

1.1 IVC from SNARK

Previously, IVC was constructed based on SNARK, and the recursive circuit included computation logic and SNARK verification logic.

Since SNARK verification logic needs to be expressed in the circuit, and SNARK verification logic often involves some circuit-unfriendly operations, such as pairing and non-native operations, the SNARK-based IVC construction has a large recursive overhead.

1.2 IVC from Folding

The IVC construction based on Folding replaces the SNARK verifier circuit with the Folding verifier circuit, greatly reducing the recursive overhead. In Nova, the prover only needs to perform 2 MSM operations of $\mathcal{O}(C)$ scale and a small amount of hashing for each step, without requiring FFT. Practice shows that Nova’s recursive overhead is around 10,000 constraints.

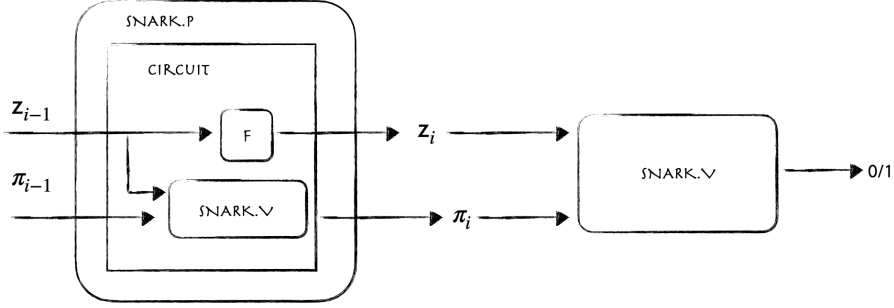


Figure 1: IVC from SNARK

2 NIVC

IVC requires the computation logic of each iteration to be the same, which has a gap with practical applications. For example, a CPU executes arbitrary instructions from the instruction set in each iteration. This requires the use of NIVC. NIVC (Non-uniform IVC) allows the computation function executed in each iteration to vary within the function set.

To construct NIVC, a naive method is to "use a universal circuit to express multiple step functions", which is the commonly used Selector scheme nowadays: laying out all instruction circuits and then activating one of the step instruction circuits through a Selector. This approach has a major drawback: the scale of the universal circuit will expand to the sum of all instruction circuits. In a VM or LLM, an instruction circuit represents a supported instruction. In this case, if the supported instruction set contains a large number of instructions, the final universal circuit will have a large expansion. This path is not ideal, so we need to find an alternative: is there a Folding scheme that supports multiple step functions??

2.1 SuperNova

The answer is SuperNova[2], an extension of Nova[1]. Its outstanding feature is that the cost of each step of the prover is only proportional to the size of the instruction circuit invoked by the program, which is extremely advantageous when the instruction set is complex.

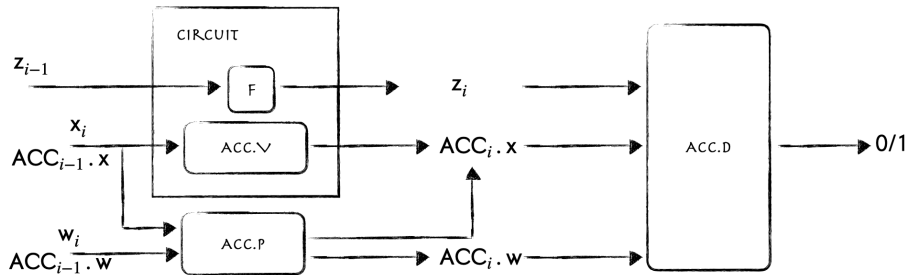


Figure 2: IVC from Folding

2.2 Folding More

SuperNova supports only R1CS. In practical applications, circuits with richer expressiveness may be needed, such as custom gates and lookups. There are many existing schemes exploring how to fold circuits with richer expressiveness, such as Sangria[4], Origami[3], HyperNova[6], and Protostar[5].

3 zkLLM

3.1 Design

We construct a RAM machine that supports an instruction set $IS = F_{i \geq \ell}$ of size ℓ . This RAM machine has 3 control registers `ts`, `pc`, `ag` for flow control; 1 general-purpose register `[gpr]` for recording the commitment of a general-purpose register of length k ; 1 stack pointer register and 1 stack register `sp`, `[sm]` for recording the commitment of a stack pointer and stack memory of length m , respectively. `in` and `out` are used to record the input and output, and `[p]` is used to record the input program.

<code>ts</code>	<code>pc</code>	<code>ag</code>	<code>[gpr]</code>	<code>sp</code>	<code>[sm]</code>	<code>in</code>	<code>out</code>	<code>[p]</code>
-----------------	-----------------	-----------------	--------------------	-----------------	-------------------	-----------------	------------------	------------------

Table 1: Registers for flow control, numerical calculations, and input and output

The design distinguishes between general-purpose registers and stack memory, referring to the design of registers and memory in typical computer architectures. When operating only on registers, the overhead is smaller because the length of registers is limited and the cost of opening their commitments is small.

For each opcode, a corresponding augmented circuit is defined:

$$F_{j \geq [\ell]}^o(\text{vk}, \mathbf{U}_i, \mathbf{u}_i, (\text{ts}_i, \text{pc}_i, \text{ip}_i, \text{ag}_i, \text{sp}_i, [\text{sm}]_i, \text{in}_0, \text{out}_i, [\text{gpr}]_i, [\mathbf{p}]_0), \bar{T}) \rightarrow x:$$

1. Update $\text{ts}_{i+1} \leftarrow \text{ts}_i + 1$
2. Calculate the instruction $[\text{ip}]_{i+1} \in Z_{\ell+1} \leftarrow \varphi(\text{pc}_i, [\mathbf{p}]_0)$ to be executed in the current step based on the pc_i input from the previous step
3. If $\text{ts}_i == 0$ (for INIT):
 - (a) Initialize the running instance list $\mathbf{U}_{i+1} \leftarrow \mathbf{u}_i^\ell$
 - (b) Check that in_0 is correctly uploaded to gpr , which is an index-lookup check, as shown in Figure 3
 - (c) Verify $[\text{sm}]_i = [0^m]$, $\text{pc} = 0$, $\text{ag} = 0$ and $\text{sp} = 0$ to ensure correct initialization
4. Otherwise:
 - (a) Verify $\mathbf{u}_i.x = \text{hash}(\text{vk}, \mathbf{U}_i, \text{ts}_i, \text{pc}_i, \text{ip}_i, \text{ag}_i, \text{sp}_i, [\text{sm}]_i, \text{in}_0, \text{out}_i, [\text{gpr}]_i, [\mathbf{p}]_0)$ to ensure the output of the previous step is the input of the current step
 - (b) Verify $j = \text{ip}_{i+1}$ to ensure the correct instruction circuit is constrained
 - (c) Verify $(\mathbf{u}_i.\bar{E}, \mathbf{u}_i.u) = (0, 1)$ to ensure the augmented circuit strictly holds
 - (d) Update the running instance list $\mathbf{U}_{i+1}[\text{ip}_i] \leftarrow \text{NIFS.V}(\text{vk}[\text{ip}_i], \mathbf{U}_i[\text{ip}_i], \mathbf{u}_i, \bar{T})$, folding the augmented circuit
5. Update the register state according to the opcode
$$(\text{pc}_{i+1}, \text{ag}_{i+1}, [\text{gpr}]_{i+1}, \text{sp}_{i+1}, [\text{sm}]_{i+1}) \leftarrow F_j(\text{pc}_i, \text{ag}_i, [\text{gpr}]_i, \text{sp}_i, [\text{sm}]_i)$$
6. Output
$$x \leftarrow \text{hash}(\text{vk}, \mathbf{U}_{i+1}, \text{ts}_{i+1}, \text{pc}_{i+1}, \text{ip}_{i+1}, \text{ag}_{i+1}, \text{sp}_{i+1}, [\text{sm}]_{i+1}, \text{in}_0, \text{out}_{i+1}, [\text{gpr}]_{i+1}, [\mathbf{p}]_0)$$

Notes:

1. $[\mathbf{p}]$, $[\text{gpr}]$, $[\text{sm}]$ represent the commitment of the input vector (general-purpose register), input code, and stack memory
2. \mathbf{p} represents the input code, which is a table consisting of $(\text{pc}, \text{ip}, \text{operand})$, as shown in Figure 4
3. The $\varphi(\cdot)$ in step 2 is a decoder that takes the program \mathbf{p} and program counter pc as input and calculates the instruction and operands to be executed in the current step, which is essentially an index-lookup, as shown in Figure 5
4. The opcode updates some registers and outputs, such as MUL and ADD modifying the pc , ip , and $[\text{gpr}]$ registers; JMP and JMPC modifying the pc and ip registers; and S_MUL modifying the sp and $[\text{sm}]$ registers

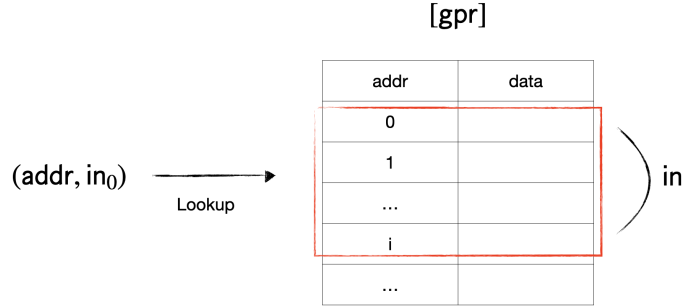


Figure 3: lookup the input in the gpr table

[p]

pc	ip	operand

Figure 4: input program table

The prover updates the proof Π using the trace at each step, $\mathcal{P}(\text{pk}, \Pi_i) \rightarrow \Pi_{i+1}$:

1. Parse the proof Π_i of step i as

$$((U_i, W_i), (u_i, w_i), (ts_i, pc_i, ip_i, ag_i, sp_i, [sm]_i, in_0, out_i, [gpr]_i, [p]_0))$$
2. If $ts_i == 0$:
 - (a) Initialize $((U_{i+1}, W_{i+1}), \bar{T}) \leftarrow (u_{\mathcal{P}}^\ell, w_{\mathcal{P}}^\ell, u_{\mathcal{P}} \cdot \bar{E})$
 - (b) Initialize $[sm]_i = [0^m]$, $pc = 0$, $ag = 0$, $sp = 0$, $[gpr]$
3. Otherwise:
 - (a) Update the corresponding running instance according to the instruction pointer

$$(U_{i+1}[ip_i], W_{i+1}[ip_i], \bar{T}) \leftarrow \text{NIFS.P}(\text{pk}[ip_i], (U_i[ip_i], W_i[ip_i]), (u_i, w_i))$$
 - (b) Calculate the instruction $[ip]_{i+1} \in Z_{\ell+1} \leftarrow \varphi(pc_i, [p]_0)$ to be executed in the current step

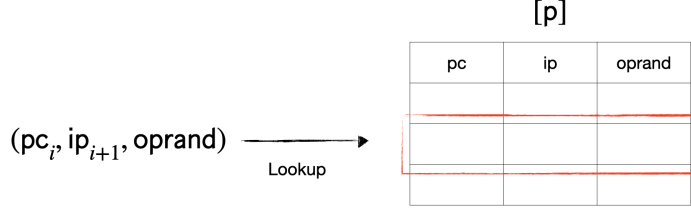


Figure 5: lookup the triple in the program table

- (c) Calculate the trace of the current step's augmented circuit
 $(u_{i+1}, w_{i+1}) \leftarrow \text{trace}(F_{ip_{i+1}}^o, vk, U_i, u_i, ts_i, pc_i, ip_i, ag_i, sp_i, [sm]_i, in_0, out_i, [gpr]_i, [p]_0, \bar{T})$
- (d) Update the proof $\Pi_{i+1} \leftarrow ((u_i, w_i), ip_{i+1})$

The verifier obtains the proof Π_i from the last folding $\mathcal{V}(vk, \Pi_i) \rightarrow \{0, 1\}$:

1. If $ts_i == 0$:
 - (a) Verify that in_0 is correctly uploaded to gpr
 - (b) Verify $[sm]_i = [0^m]$, $pc = 0$, $ag = 0$, $sp = 0$ to ensure correct initial values
2. Otherwise:
 - (a) Parse the proof Π_i as
 $((U_i, W_i), (u_i, w_i), (ts_i, pc_i, ip_i, ag_i, sp_i, [sm]_i, in_0, out_i, [gpr]_i, [p]_0))$
 - (b) Verify $u_i.x = \text{hash}(vk, U_i, ts_i, pc_i, ip_i, ag_i, sp_i, [sm]_i, in_0, out_i, [gpr]_i, [p]_0)$
 - (c) $\varphi(pc_i, [p]_0) = \text{endvar}$
 - (d) Verify $(u_i.\bar{E}, u_i.u) = (0, 1)$
 - (e) For $F_{ip_i}^o$, verify the (u_i, w_i) of the last iteration
 - (f) For all F^o , verify all $(U_i[j], W_j)_{j \geq [\ell]}$

3.2 Opcode Examples

3.2.1 INIT

INIT is used to upload the initial input in_0 agreed upon by both parties (including the consensus of private input) to the predetermined addresses of gpr . The INIT circuit needs to constrain:

1. $ts_i = 0$, $[sm]_i = [0^m]$, $pc = 0$, $ag = 0$, $sp = 0$
2. The values at the predetermined positions of gpr are equal to the input, $in_0 = OPEN(gpr_i, addr)$, where $OPEN$ can be implemented as the index lookup constraint shown in Figure 3
3. Update $pc_{i+1} = pc_i + 1$, and the rest of the register states remain unchanged

Note: The maximum length of the initial input is specified as d , and less than that is padded with 0. For example, if the maximum initial length is $d = 4$, and the initial input is $in_0 = [1, 2, 0, 0]$, then 2 zeros are padded. The initial register state satisfying the constraints should be:

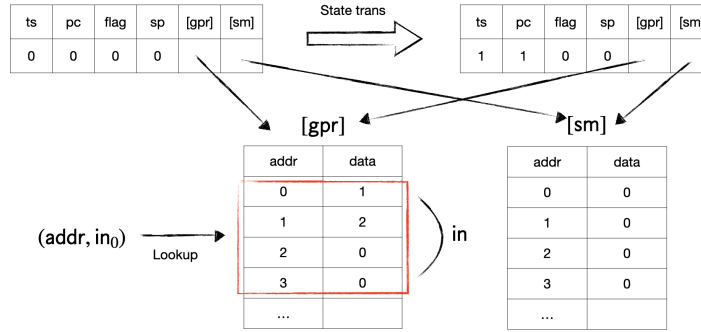


Figure 6: INIT legal state

3.2.2 ADD1_4

ADD1.4 $addr_0 \ addr_1 \ addr_2$ is used for the addition of two $1*4$ tensors at specified addresses, and the result is stored at the specified address. The update requires constraints:

1. The computation at the specified addresses of gpr is correct: $left \leftarrow OPEN([gpr]_i, addr_0)$, $right \leftarrow OPEN([gpr]_i, addr_1)$, $output \leftarrow OPEN([gpr]_{i+1}, addr_2)$, $output = left + right$
2. Update $pc_{i+1} = pc_i + 1$, $[gpr]_{i+1} = UPDATE([gpr]_i, addr_2)$

Note: \leftarrow represents generating an intermediate variable and constraining it. The $UPDATE$ constraint is $[gpr]_{i+1} = [gpr]_i + \sum addr_2([gpr]_{i+1}^{addr_2} - [gpr]_i^{addr_2}) \cdot [L^{addr_2}(X)]$, and its complexity is related to the number of modified addresses.

For example, $ADD1.4 \ 0 \ 4 \ 8$ adds $[1, 2, 0, 0]$ at address 0 and $[1, 2, 3, 4]$ at address 4, and places the result at address 8.

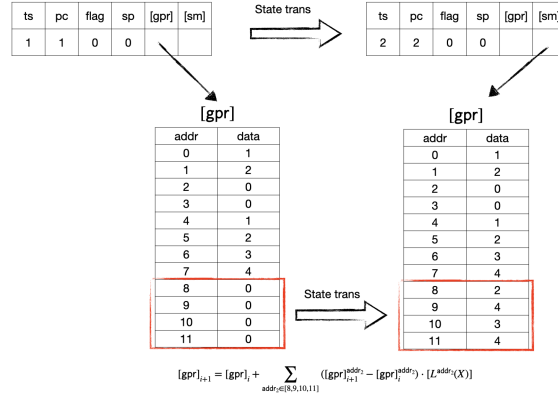


Figure 7: ADD1.4 legal state (ignoring the unchanged stack memory and output)

3.2.3 LE1.4

LE1.4 $addr_0$ $addr_1$ is used to compare two 1×4 tensors at specified addresses, and the flag is updated based on the comparison result. The circuit needs to constrain:

1. If $a \geq b$: $ag_{i+1} = 1$
2. Update $pc_{i+1} = pc_i + 1$

For example, LE1.4 0 4 compares $[1, 2, 0, 0]$ at address 0 with $[1, 2, 3, 4]$ at address 4.

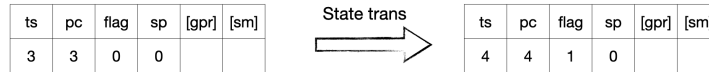


Figure 8: LE1.4 legal state (ignoring the unchanged registers, stack memory, and output)

3.2.4 JMP

JMPC $addr_0$ $addr_1$ is used to jump based on ag . The circuit needs to constrain:

1. If ag_{i+1} holds, update $pc_{i+1} = addr_0$, otherwise update $pc_{i+1} = addr_1$

For example, JMPC 1 4 sets pc to 1.

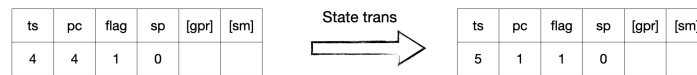


Figure 9: JMPC legal state (ignoring the unchanged registers, stack memory, and output)

3.2.5 RETURN

RETURN $addr$ is used to download the data at the specified address in $addr$ to out . The circuit needs to constrain:

1. Update $out_{i+1} = OPEN([gpr]_i, addr)$
2. Update $pc_{i+1} = pc_{end}$

The maximum length d of the output out is specified. For example, RETURN 8 outputs $[2, 4, 3, 4]$ at address 8 to out .

3.2.6 END

END is used to indicate the end of the program execution. The circuit needs to constrain:

1. All states remain unchanged

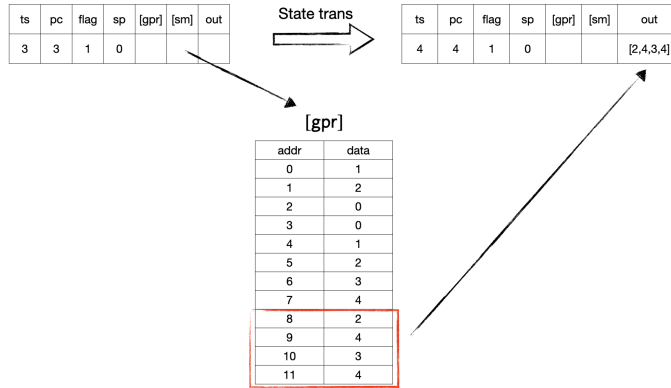


Figure 10: RETURN legal state (ignoring the unchanged stack memory)

References

- [1] Benedikt Bünz, Yuncong Hu, and Shin'ichiro Matsuo. Nova: Recursive zero-knowledge arguments from folding schemes. Cryptology ePrint Archive, Paper 2021/370, 2021. <https://eprint.iacr.org/2021/370>.
- [2] Benedikt Bünz, Yuncong Hu, and Shin'ichiro Matsuo. Supernova: Snarks for nive with linear complexity overhead. Cryptology ePrint Archive, Paper 2022/1758, 2022. <https://eprint.iacr.org/2022/1758>.
- [3] Geometry Labs. Origami, 2023.
- [4] Geometry Labs. Sangria: A folding scheme for plonk, 2023.
- [5] Chuanjiang Li and Hongzhao Chen. Protostar: Folding starks with linear complexity. Cryptology ePrint Archive, Paper 2023/620, 2023. <https://eprint.iacr.org/2023/620>.
- [6] Thomas Sirvent, Yuncong Hu, and Geoffroy Couteau. Hypernova: Folding hyperplonk. Cryptology ePrint Archive, Paper 2023/573, 2023. <https://eprint.iacr.org/2023/573>.