# Folding based zkLLM

# 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 $y = F^{(n)}(x)$ .

# IVC from SNARK

Previously, IVC was constructed based on SNARK, and the recursive circuit included computation logic and SNARK verification logic.
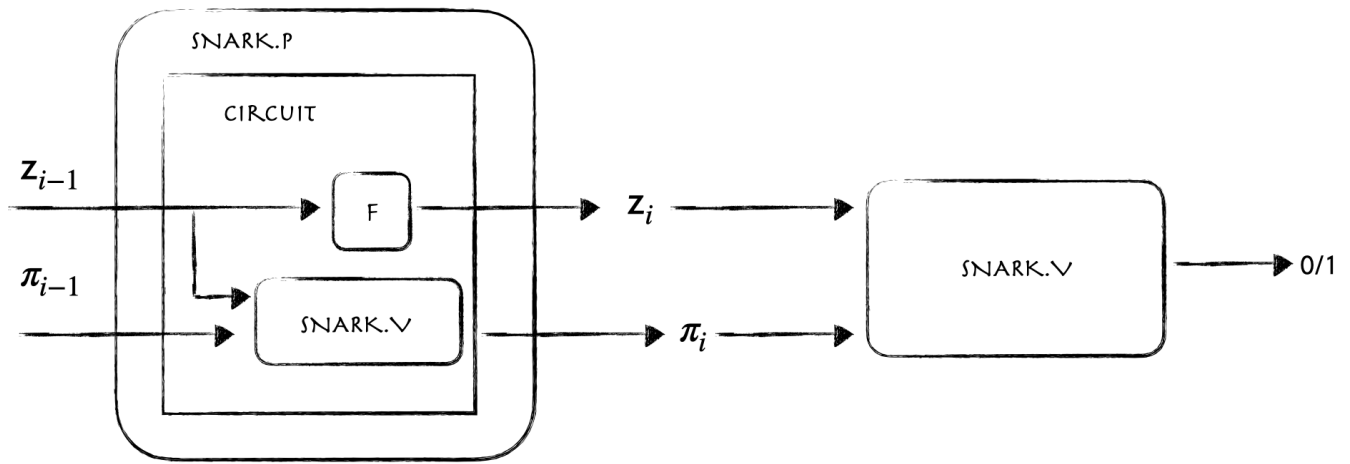
Fig 1.1

Since the 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.

# 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.
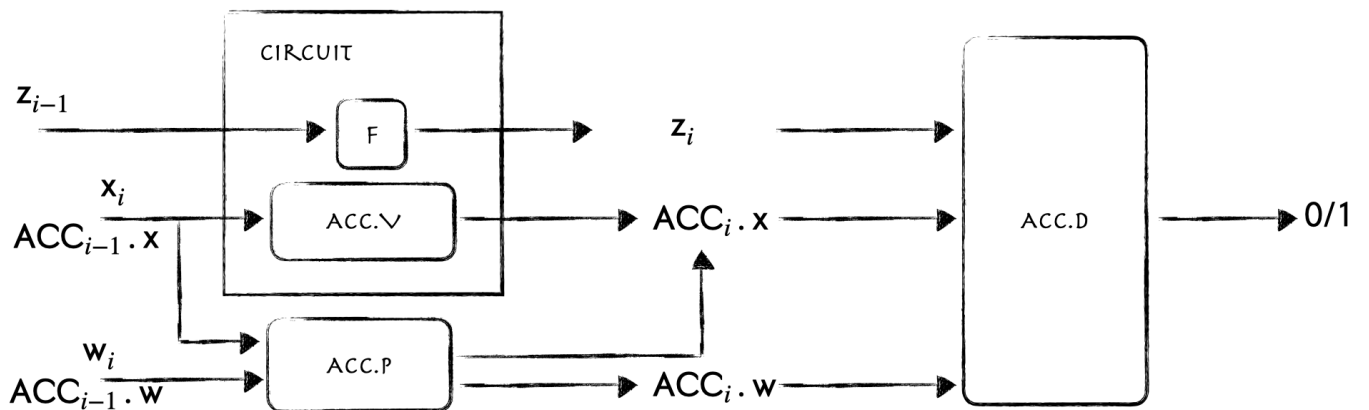


Fig 1.2

# 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 good, so we have to find another way: is there a Folding scheme that supports multiple step functions?

## SuperNova

The answer is SuperNova[1], an extension of Nova[2]. 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.

## Folding More

SuperNova only supports 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[3], Origami[4], HyperNova[5], and Protostar[6].

# zkLLM

## Design

We construct a RAM machine that supports an instruction set $\mathsf{IS} = F_{i \in [\ell]}$ of size $\ell$. This RAM machine has 3 control registers **ts**, **pc**, **flag** for flow control; 1 general–purpose

register $[\mathbf{gpr}]$ for recording the commitment of a general–purpose register of length $k$ ; 1 stack pointer register and 1 stack register $\mathbf{sp}$ , $[\mathbf{sm}]$ for recording the commitment of a stack pointer and stack memory of length $m$ , respectively. $\mathbf{in}$ and $\mathbf{out}$ are used to record input and output, and $[\mathbf{p}]$ is used to record the input program.

| ts | pc | flag | [gpr] | sp | [sm] | in | out | [p] |
|----|----|------|-------|----|------|----|-----|-----|
|    |    |      |       |    |      |    |     |     |

<div align="center">Table 3.1</div>

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, an corresponding augmented circuit is defined:

$$F'_{j\in[\ell]}(\mathsf{vk}, \mathsf{U}_i, \mathsf{u}_i, (\mathbf{ts}_i, \mathbf{pc}_i, \mathbf{ip}_i, \mathbf{flag}_i, \mathbf{sp}_i, [\mathbf{sm}]_i, \mathbf{in}_0, \mathbf{out}_i, [\mathbf{gpr}]_i, [\mathbf{p}]_0), \bar{T}) \to x :$$

1. Update $\mathbf{ts}_{i+1} \leftarrow \mathbf{ts}_i + 1$

2. Calculate the instruction $[\mathbf{ip}]_{i+1} \in \mathbb{Z}^*_{\ell+1} \leftarrow \varphi(\mathbf{pc}_i, [\mathbf{p}]_0)$ to be executed in the current step based on the $\mathbf{pc}_i$ input from the previous step

3. If $\mathbf{ts}_i == 0$ (for INIT):

    a. Initialize the running instance list $\mathsf{U}_{i+1} \leftarrow \mathsf{u}^\ell_\perp$

    b. Check that $\mathbf{in}_0$ is correctly uploaded to $\mathbf{gpr}$ , which is a index–lookup check, as shown in Figure 3.1



<div align="center">Fig 3.1</div>

c. Verify $[\mathbf{sm}]_i = [\mathbf{0}^m]$ , $\mathbf{pc} = 0$ , $\mathbf{flag} = 0$ , $\mathbf{sp} = 0$ to ensure correct initialization

Otherwise:

4. Verify $\mathbf{u}_i.x = \mathsf{hash}(\mathbf{vk}, \mathbf{U}_i, \mathbf{ts}_i, \mathbf{pc}_i, \mathbf{ip}_i, \mathbf{flag}_i, \mathbf{sp}_i, [\mathbf{sm}]_i, \mathbf{in}_0, \mathbf{out}_i, [\mathbf{gpr}]_i, [\mathbf{p}]_0)$ to ensure the output of the previous step is the input of the current step

5. Verify $j = \mathbf{ip}_{i+1}$ to ensure the correct instruction circuit is constrained

6. Verify $(\mathbf{u}_i.\bar{E}, \mathbf{u}_i.u) = (0, 1)$ to ensure the augmented circuit strictly holds

7. Update the running instance list $\mathbf{U}_{i+1}[\mathbf{ip}_i] \leftarrow \mathsf{NIFS.V}(\mathbf{vk}[\mathbf{ip}_i], \mathbf{U}_i[\mathbf{ip}_i], \mathbf{u}_i, \bar{T})$ , folding the augmented circuit

8. Update the register state according to the opcode
$(\mathbf{pc}_{i+1}, \mathbf{flag}_{i+1}, [\mathbf{gpr}]_{i+1}, \mathbf{sp}_{i+1}, [\mathbf{sm}]_{i+1}) \leftarrow F_j(\mathbf{pc}_i, \mathbf{flag}_i, [\mathbf{gpr}]_i, \mathbf{sp}_i, [\mathbf{sm}]_i)$

9. Output
$x \leftarrow \mathsf{hash}(\mathbf{vk}, \mathbf{U}_{i+1}, \mathbf{ts}_{i+1}, \mathbf{pc}_{i+1}, \mathbf{ip}_{i+1}, \mathbf{flag}_{i+1}, \mathbf{sp}_{i+1}, [\mathbf{sm}]_{i+1}, \mathbf{in}_0, \mathbf{out}_{i+1}, [\mathbf{gpr}]_{i+1}, [\mathbf{p}]_0)$

Notes:

1. $[\mathbf{p}]$ , $[\mathbf{gpr}]$ , $[\mathbf{sm}]$ 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 $(\mathbf{pc}, \mathbf{ip}, \mathbf{oprand})$ , as shown in Figure 3.2

# [p]

| pc | ip | oprand |
|----|----|--------|
|    |    |        |
|    |    |        |
|    |    |        |

Fig 3.2

3. The $\varphi(\cdot)$ in step 2 is a decoder that takes the program $\mathbf{p}$ and program counter $\mathbf{pc}$ as input and calculates the instruction and operands to be executed in the current step, which is essentially an index–lookup
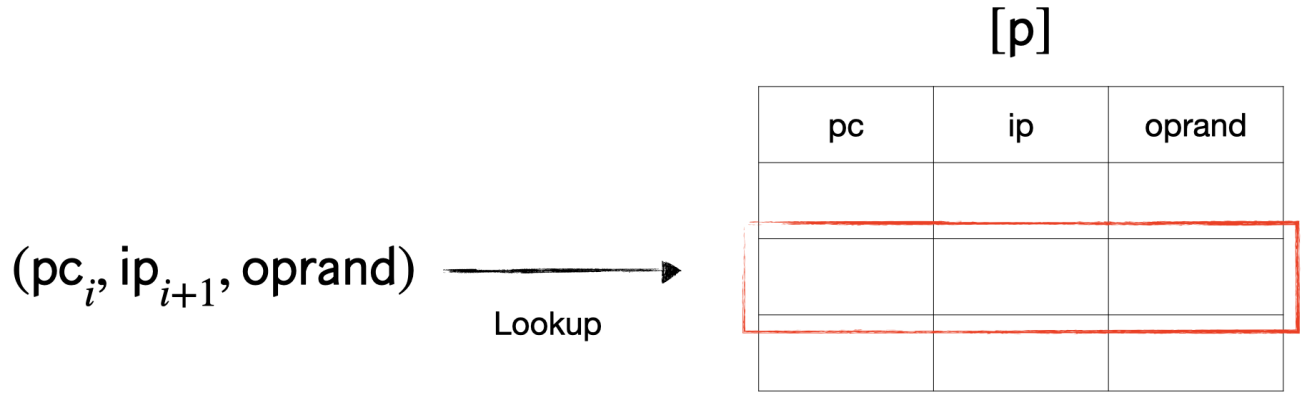
**[p]**

| | pc | ip | oprand |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |

$$(\mathsf{pc}_i, \mathsf{ip}_{i+1}, \mathsf{oprand}) \xrightarrow{\quad \text{Lookup} \quad}$$

Fig 3.3

4. The opcode updates some registers and outputs, such as **MUL** and **ADD** modifying the **pc** **ip** $[\mathsf{gpr}]$ registers, **JMP** and **JMPC** modifying the **pc** **ip** registers, and **S_MUL** modifying the **sp** $[\mathsf{sm}]$ registers

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

1. Parse the proof $\Pi_i$ of step $i$ as
   $$((\mathsf{U}_i, \mathsf{W}_i), (\mathsf{u}_i, \mathsf{w}_i), (\mathsf{ts}_i, \mathsf{pc}_i, \mathsf{ip}_i, \mathsf{flag}_i, \mathsf{sp}_i, [\mathsf{sm}]_i, \mathsf{in}_0, \mathsf{out}_i, [\mathsf{gpr}]_i, [\mathsf{p}]_0))$$

2. If $\mathsf{ts}_i == 0$ :

   a. Initialize $((\mathsf{U}_{i+1}, \mathsf{W}_{i+1}), \bar{T}) \leftarrow (\mathsf{u}_\perp^\ell, \mathsf{w}_\perp^\ell, \mathsf{u}_\perp.\bar{E})$

   b. Initialize $[\mathsf{sm}]_i = [0^m]$ , $\mathsf{pc} = 0$ , $\mathsf{flag} = 0$ , $\mathsf{sp} = 0$ , $[\mathsf{gpr}]$

Otherwise:

3. Update the corresponding running instance according to the instruction pointer
   $$(\mathsf{U}_{i+1}[\mathsf{ip}_i], \mathsf{W}_{i+1}[\mathsf{ip}_i], \bar{T}) \leftarrow \mathsf{NIFS.P}(\mathsf{pk}[\mathsf{ip}_i], (\mathsf{U}_i[\mathsf{ip}_i], \mathsf{W}_i[\mathsf{ip}_i]), (\mathsf{u}_i, \mathsf{w}_i))$$

4. Calculate the instruction $[\mathsf{ip}]_{i+1} \in \mathbb{Z}_{\ell+1}^* \leftarrow \varphi(\mathsf{pc}_i, [\mathsf{p}]_0)$ to be executed in the current step

5. Calculate the trace of the current step's augmented circuit
   $$(\mathsf{u}_{i+1}, \mathsf{w}_{i+1}) \leftarrow$$
   $$\mathsf{trace}(F'_{\mathsf{ip}_{i+1}}, \mathsf{vk}, \mathsf{U}_i, \mathsf{u}_i, \mathsf{ts}_i, \mathsf{pc}_i, \mathsf{ip}_i, \mathsf{flag}_i, \mathsf{sp}_i, [\mathsf{sm}]_i, \mathsf{in}_0, \mathsf{out}_i, [\mathsf{gpr}]_i, [\mathsf{p}]_0, \bar{T})$$

6. Update the proof $\Pi_{i+1} \leftarrow ((\mathsf{u}_i, \mathsf{w}_i), \mathsf{ip}_{i+1})$

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

If $\mathbf{ts}_i == 0$ :

1. Verify that $\mathbf{in}_0$ is correctly uploaded to $\mathbf{gpr}$

2. Verify $[\mathbf{sm}]_i = [0^m]$ , $\mathbf{pc} = 0$ , $\mathbf{flag} = 0$ , $\mathbf{sp} = 0$ to ensure correct initial values

Otherwise:

3. Parse the proof $\Pi_i$ as
$$((\mathbf{U}_i, \mathbf{W}_i), (\mathbf{u}_i, \mathbf{w}_i), (\mathbf{ts}_i, \mathbf{pc}_i, \mathbf{ip}_i, \mathbf{flag}_i, \mathbf{sp}_i, [\mathbf{sm}]_i, \mathbf{in}_0, \mathbf{out}_i, [\mathbf{gpr}]_i, [\mathbf{p}]_0))$$

4. Verify $\mathbf{u}_i.x = \mathsf{hash}(\mathsf{vk}, \mathbf{U}_i, \mathbf{ts}_i, \mathbf{pc}_i, \mathbf{ip}_i, \mathbf{flag}_i, \mathbf{sp}_i, [\mathbf{sm}]_i, \mathbf{in}_0, \mathbf{out}_i, [\mathbf{gpr}]_i, [\mathbf{p}]_0)$

5. $\varphi(\mathbf{pc}_i, [\mathbf{p}]_0) = \mathsf{endvar}$

6. Verify $(\mathbf{u}_i.\bar{E}, \mathbf{u}_i.u) = (0, 1)$

7. For $F'_{\mathsf{ip}_i}$ , verify the $(\mathbf{u}_i, \mathbf{w}_i)$ of the last iteration

8. For all $F'$ , verify all $(\mathbf{U}_i[j], \mathbf{W}_j)_{j \in [\ell]}$

# Opcode Examples

## INIT

INIT is used to upload the initial input $\mathbf{in}_0$ agreed upon by both parties (including the consensus of private input) to the predetermined addresses of $\mathbf{gpr}$ . The $\mathsf{INIT}$ circuit needs to constrain:

1. $\mathbf{ts}_i = 0$ , $[\mathbf{sm}]_i = [0^m]$ , $\mathbf{pc} = 0$ , $\mathbf{flag} = 0$ , $\mathbf{sp} = 0$

2. The values at the predetermined positions of $\mathbf{gpr}$ are equal to the input,
$\mathbf{in}_0 = \mathsf{OPEN}(\mathbf{gpr}_i, \mathbf{addr})$ , where $\mathsf{OPEN}$ can be implemented as the index lookup constraint shown in Figure 3.1

3. Update $\mathbf{pc}_{i+1} = \mathbf{pc}_i + 1$ , 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 $\mathbf{in}_0 = [1, 2, 0, 0]$ , then 2 zeros are padded. The initial register state satisfying the constraints should be:
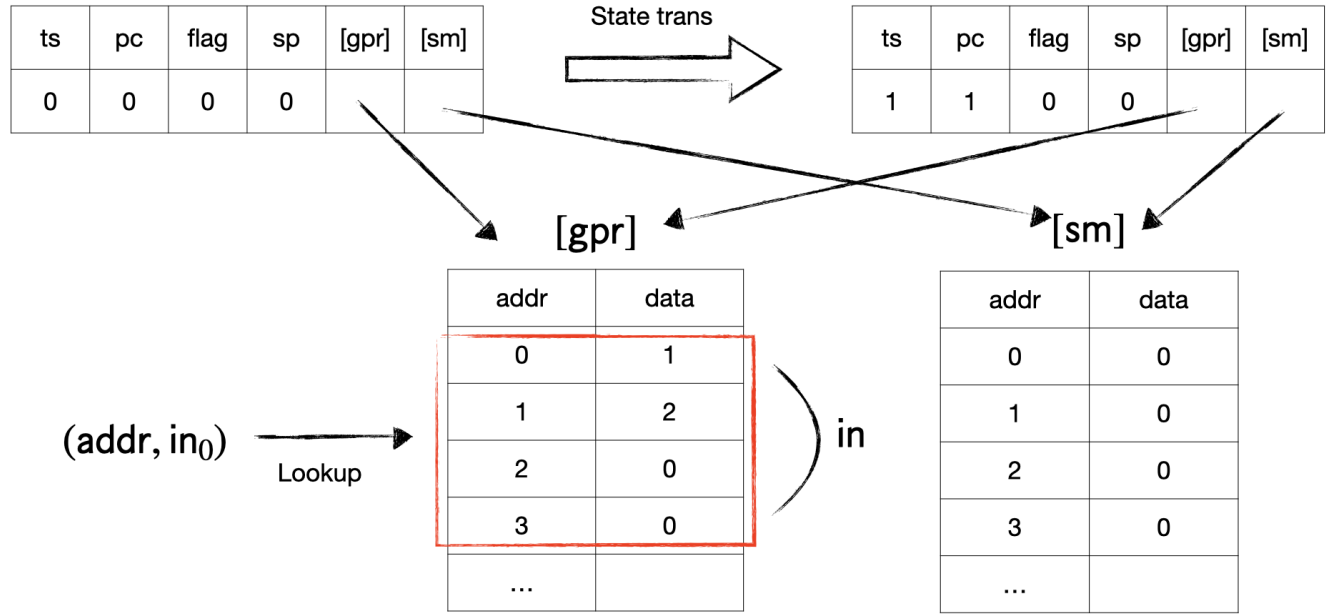
| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|------|----|-------|------|
| 0 | 0 | 0 | 0 | | |

State trans

| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|------|----|-------|------|
| 1 | 1 | 0 | 0 | | |

**[gpr]**

| addr | data |
|------|------|
| 0 | 1 |
| 1 | 2 |
| 2 | 0 |
| 3 | 0 |
| ... | |

**[sm]**

| addr | data |
|------|------|
| 0 | 0 |
| 1 | 0 |
| 2 | 0 |
| 3 | 0 |
| ... | |

$(\text{addr}, \text{in}_0)$ →  Lookup

in

Fig 3.4: INIT legal state

## ADD1_4

**ADD1_4** $\mathbf{addr}_0$ $\mathbf{addr}_1$ $\mathbf{addr}_2$ 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:
   $\mathbf{left} \leftarrow \text{OPEN}([\mathbf{gpr}]_i, \mathbf{addr}_0)$ , $\mathbf{right} \leftarrow \text{OPEN}([\mathbf{gpr}]_i, \mathbf{addr}_1)$ ,
   $\mathbf{output} \leftarrow \text{OPEN}([\mathbf{gpr}]_{i+1}, \mathbf{addr}_2)$ , $\mathbf{output} = \mathbf{left} + \mathbf{right}$

2. Update $\mathbf{pc}_{i+1} = \mathbf{pc}_i + 1$ , $[\mathbf{gpr}]_{i+1} = \text{UPDATE}([\mathbf{gpr}]_i, \mathbf{addr}_2)$

Note: $\leftarrow$ represents generating an intermediate variable and constraining it.

The **UPDATE** constraint is $[\mathbf{gpr}]_{i+1} = [\mathbf{gpr}]_i + \sum_{\mathbf{addr}_2}([\mathbf{gpr}]_{i+1}^{\mathbf{addr}_2} - [\mathbf{gpr}]_i^{\mathbf{addr}_2}) \cdot [L^{\mathbf{addr}_2}(X)]$ ,

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.
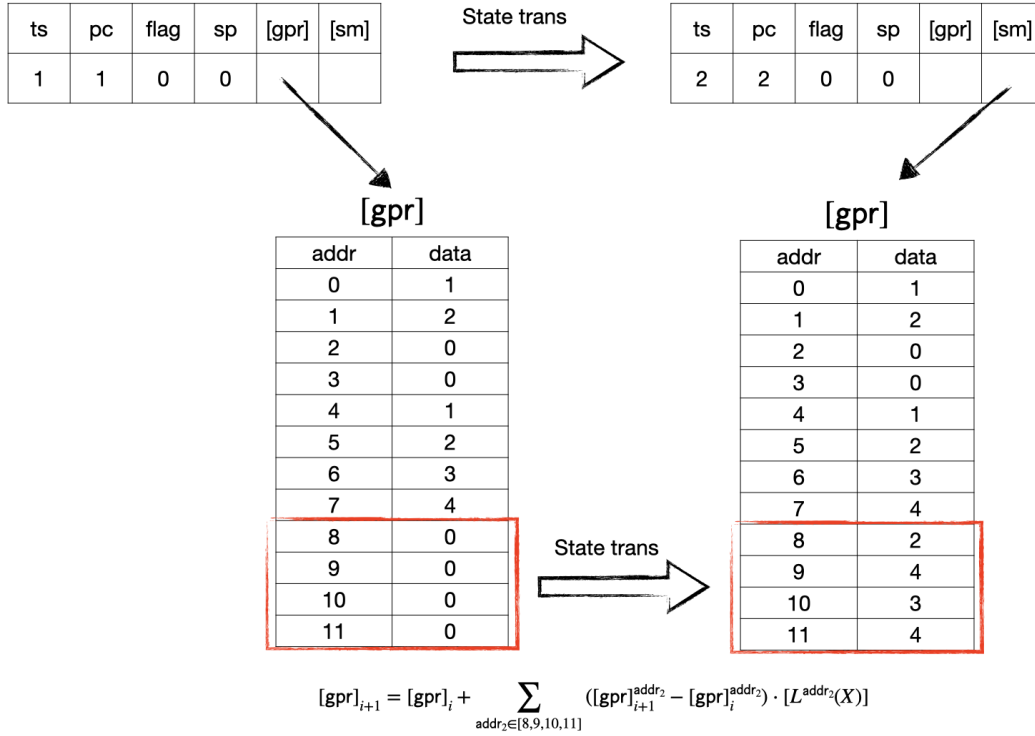
| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|----|----|----|----|
| 1 | 1 | 0 | 0 | | |

State trans →

| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|----|----|----|----|
| 2 | 2 | 0 | 0 | | |

[gpr]

| addr | data |
|----|----|
| 0 | 1 |
| 1 | 2 |
| 2 | 0 |
| 3 | 0 |
| 4 | 1 |
| 5 | 2 |
| 6 | 3 |
| 7 | 4 |
| 8 | 0 |
| 9 | 0 |
| 10 | 0 |
| 11 | 0 |

State trans →

[gpr]

| addr | data |
|----|----|
| 0 | 1 |
| 1 | 2 |
| 2 | 0 |
| 3 | 0 |
| 4 | 1 |
| 5 | 2 |
| 6 | 3 |
| 7 | 4 |
| 8 | 2 |
| 9 | 4 |
| 10 | 3 |
| 11 | 4 |

$$[\text{gpr}]_{i+1} = [\text{gpr}]_i + \sum_{\text{addr}_2 \in [8,9,10,11]} ([\text{gpr}]_{i+1}^{\text{addr}_2} - [\text{gpr}]_i^{\text{addr}_2}) \cdot [L^{\text{addr}_2}(X)]$$

Fig 3.5: ADD1_4 legal state (ignoring the unchanged stack memory and output)

## LE1_4

LE1_4 $\textbf{addr}_0$ $\textbf{addr}_1$ 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 ⩾ b: $\textbf{flag}_{i+1} = 1$
2. Update $\textbf{pc}_{i+1} = \textbf{pc}_i + 1$

For example, LE1_4 0 4 compares [1,2,0,0] at address 0 with [1,2,3,4] at address 4.
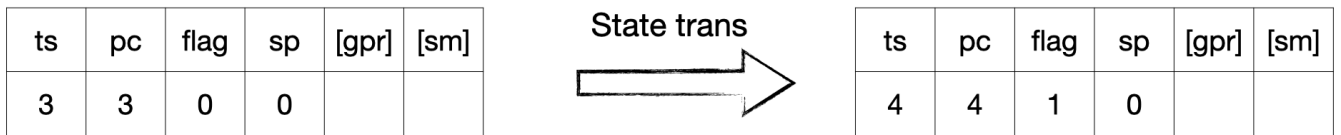
| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|----|----|----|----|
| 3 | 3 | 0 | 0 | | |

State trans →

| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|----|----|----|----|
| 4 | 4 | 1 | 0 | | |

Fig 3.6: LE1_4 legal state (ignoring the unchanged registers, stack memory, and output)

## JMPC

JMPC $\textbf{addr}_0$ $\textbf{addr}_1$ is used to jump based on flag . The circuit needs to constrain:

1. If $\textbf{flag}_{i+1}$ holds, update $\textbf{pc}_{i+1} = \textbf{addr}_0$ , otherwise update $\textbf{pc}_{i+1} = \textbf{addr}_1$

For example, $\mathsf{JMPC}\,1\,4$ sets **pc** to 1

| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|------|----|-------|------|
| 4 | 4 | 1 | 0 | | |

**State trans**

| ts | pc | flag | sp | [gpr] | [sm] |
|----|----|------|----|-------|------|
| 5 | 1 | 1 | 0 | | |

FIg 3.7: $\mathsf{JMPC}$ legal state (ignoring the unchanged registers, stack memory, and output)

## RETURN

$\mathsf{RETURN}\,\mathbf{addr}$ is used to download the data at the specified address in **addr** to **out** . The circuit needs to constrain:

1. Update $\mathbf{out}_{i+1} = \mathrm{OPEN}([\mathbf{gpr}]_i, \mathbf{addr})$
2. Update $\mathbf{pc}_{i+1} = \mathbf{pc}_{\mathrm{end}}$

The maximum length $d$ of the output **out** is specified.

For example, $\mathsf{RETURN}\,8$ outputs $[2, 4, 3, 4]$ at address 8 to **out**

| ts | pc | flag | sp | [gpr] | [sm] | out |
|----|----|------|----|-------|------|-----|
| 3 | 3 | 1 | 0 | | | |

**State trans**

| ts | pc | flag | sp | [gpr] | [sm] | out |
|----|----|------|----|-------|------|-----|
| 4 | 4 | 1 | 0 | | | [2,4,3,4] |

**[gpr]**

| addr | data |
|------|------|
| 0 | 1 |
| 1 | 2 |
| 2 | 0 |
| 3 | 0 |
| 4 | 1 |
| 5 | 2 |
| 6 | 3 |
| 7 | 4 |
| 8 | 2 |
| 9 | 4 |
| 10 | 3 |
| 11 | 4 |

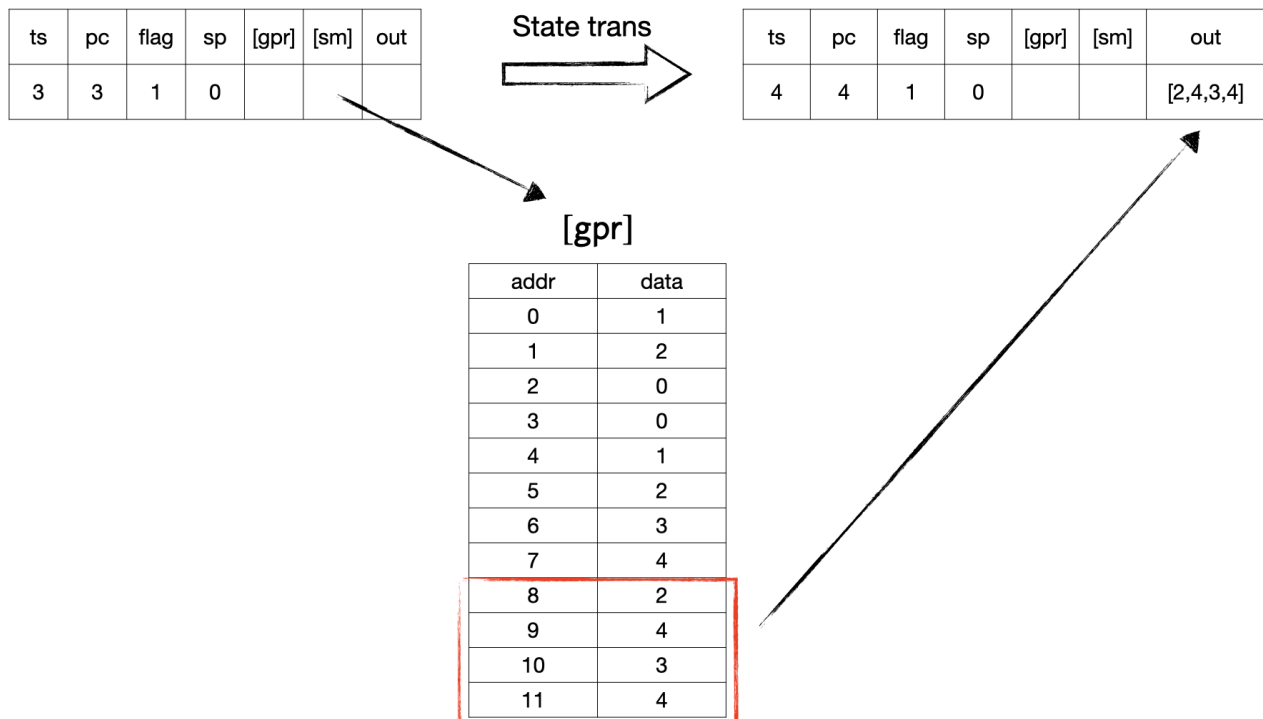Fig 3.8: $\mathsf{RETURN}$ legal state (ignoring the unchanged stack memory)

## END

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

1. All states remain unchanged

# References

[1] https://eprint.iacr.org/2022/1758.pdf

[2] https://eprint.iacr.org/2022/1758.pdf

[3] https://geometry.xyz/notebook/sangria-a-folding-scheme-for-plonk

[4] https://hackmd.io/@aardvark/rkHqa3NZ2

[5] https://eprint.iacr.org/2023/573.pdf

[6] https://eprint.iacr.org/2023/620.pdf