



Sangria is relaxed PLONK:
A Nova-style folding scheme for PLONK

APR | 2023



Sangria: Overview

- Folding schemes were introduced in [Nova](#) and are the key ingredient to achieving cheap recursion. Nova only shows a folding scheme for R1CS.
- A folding scheme compresses two instances of a problem into a single instance of the same problem
- Sangria is a folding scheme for PLONKish circuits
- Costs are similar to Nova:
 - Verifier work is constant in the depth of the circuit
 - Constants are worse than Nova, this is the price we pay to get a more flexible arithmetization



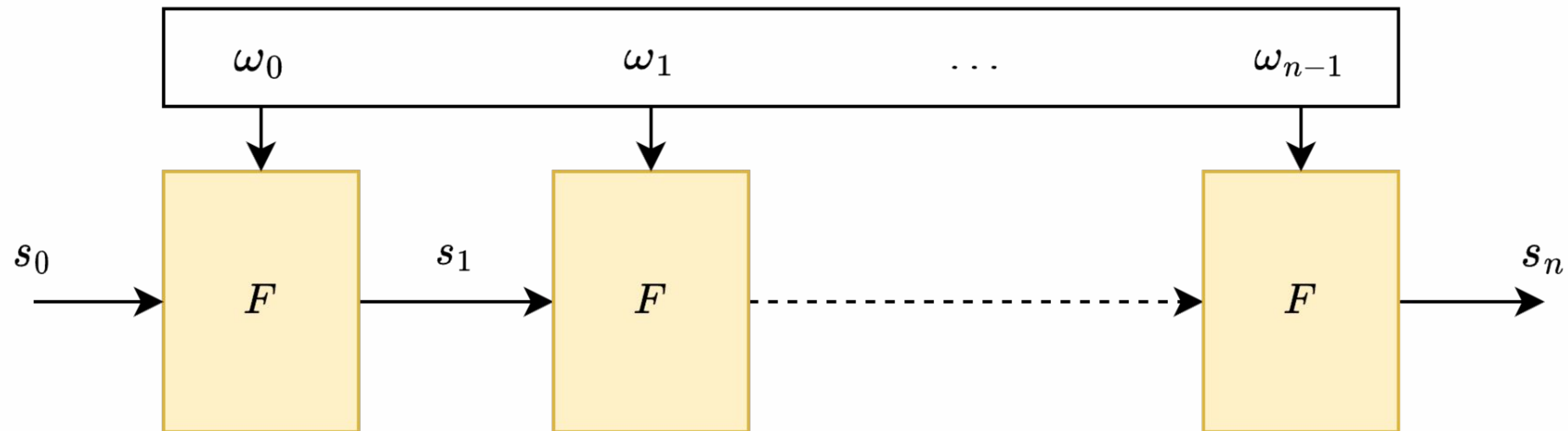
This Talk

- Motivation: why do we care about recursion? What are the known results?
- Some background: PLONK circuits and folding schemes
- Sangria: how to fold PLONK circuits
- Trade-offs and a new design space



Goal: Incrementally Verifiable Computation (IVC)

Imagine a long computation that results from n applications of a function F .

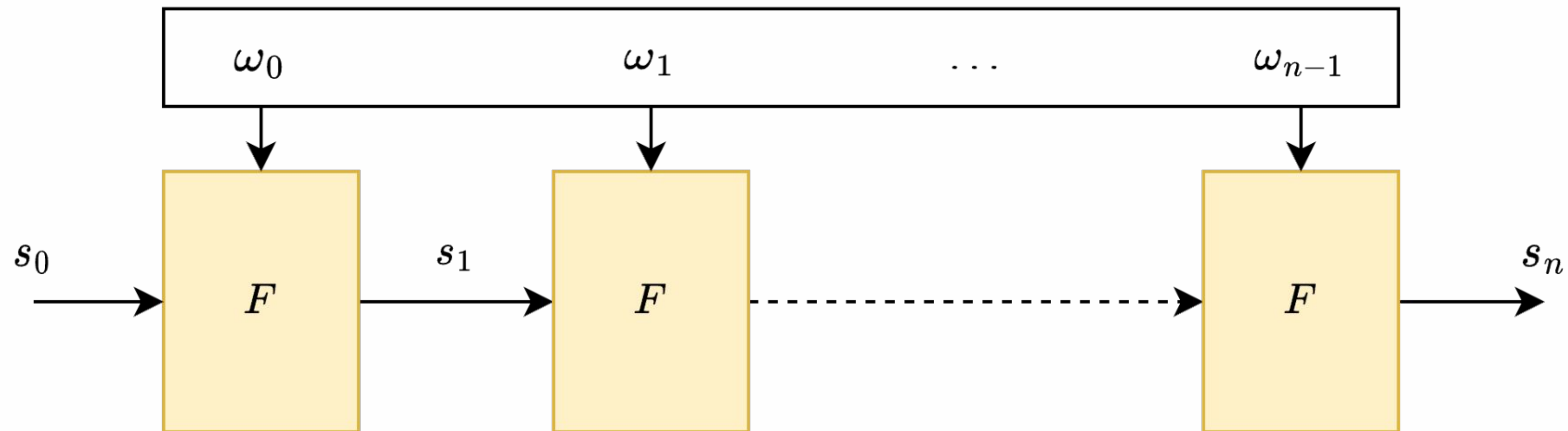


We want to produce a succinct proof that the prover knows $\omega_0, \omega_1, \dots, \omega_{n-1}$ such that s_n is the correct final state



Goal: Incrementally Verifiable Computation (IVC)

Imagine a long computation that results from n applications of a function F .



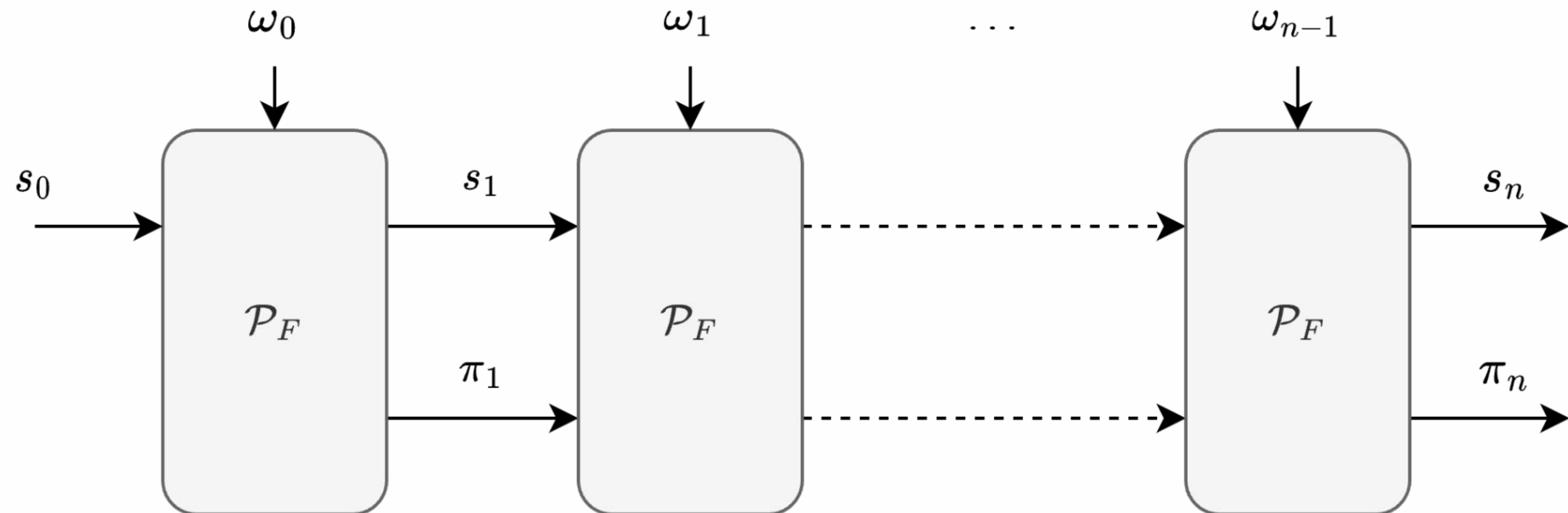
We want to produce a succinct proof that the prover knows $\omega_0, \omega_1, \dots, \omega_{n-1}$ such that s_n is the correct final state

In practice, IVC can be used to build zkVMs, rollups and verifiable delay functions



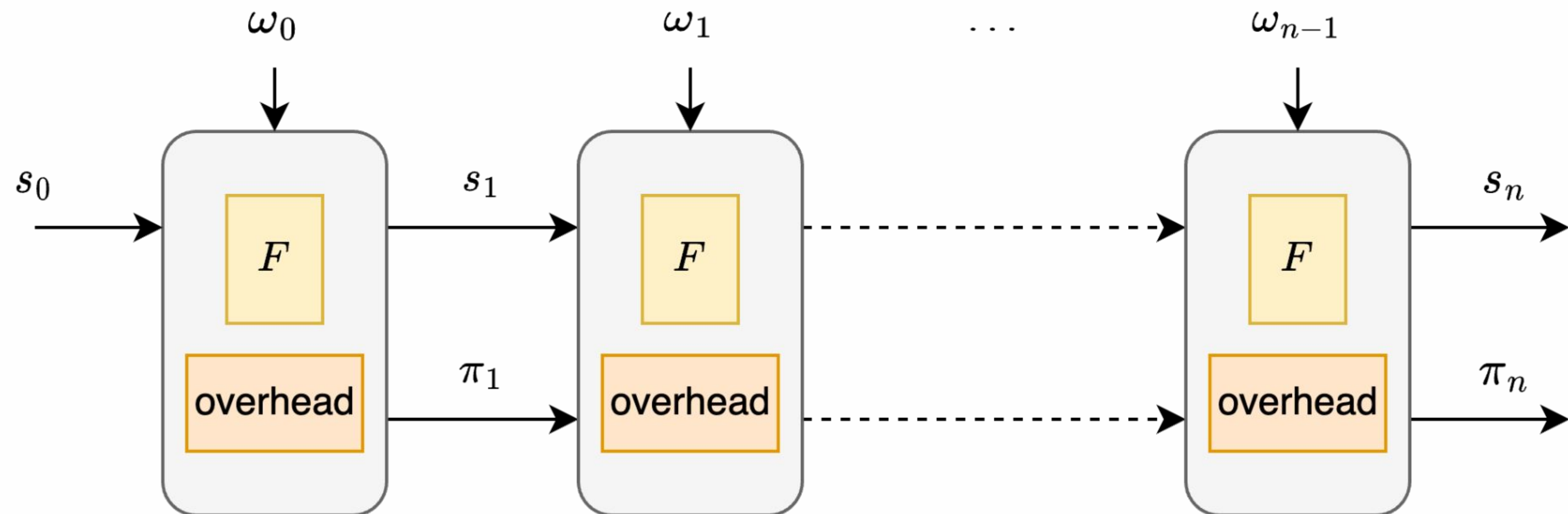
Realising IVC

To realise IVC, we define a prover that takes in a state **and** proof, and updates both



Realising IVC

To realise IVC, we define a prover that takes in a state **and** proof, and updates both

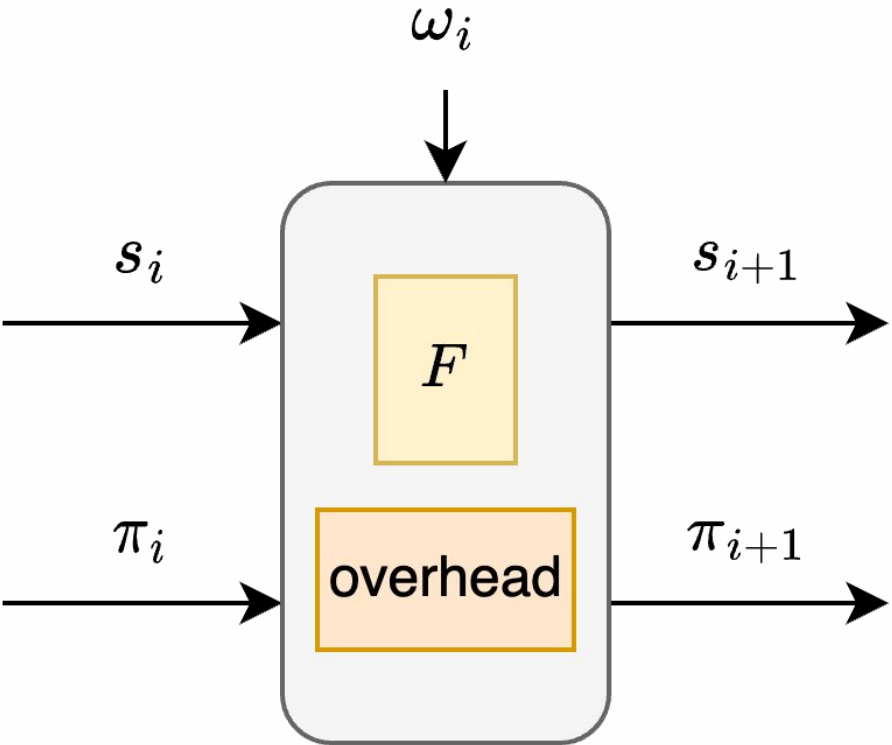


We must pay some overhead to update the proof. How small can we get it?



Recursion Overheads

Recent works have drastically reduced the recursion overhead:

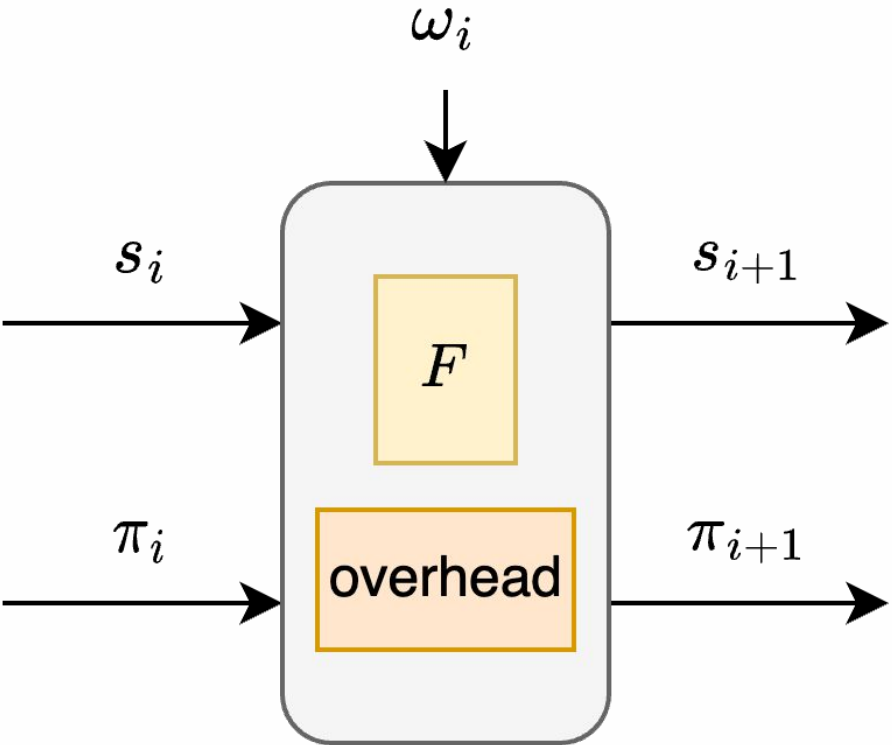


	Overhead	Examples
SNARK of SNARK	Read whole proof, run full verifier, defer nothing	Plonky2 , recursive STARK, Fractal , Groth16 + BCTV14



Recursion Overheads

Recent works have drastically reduced the recursion overhead:

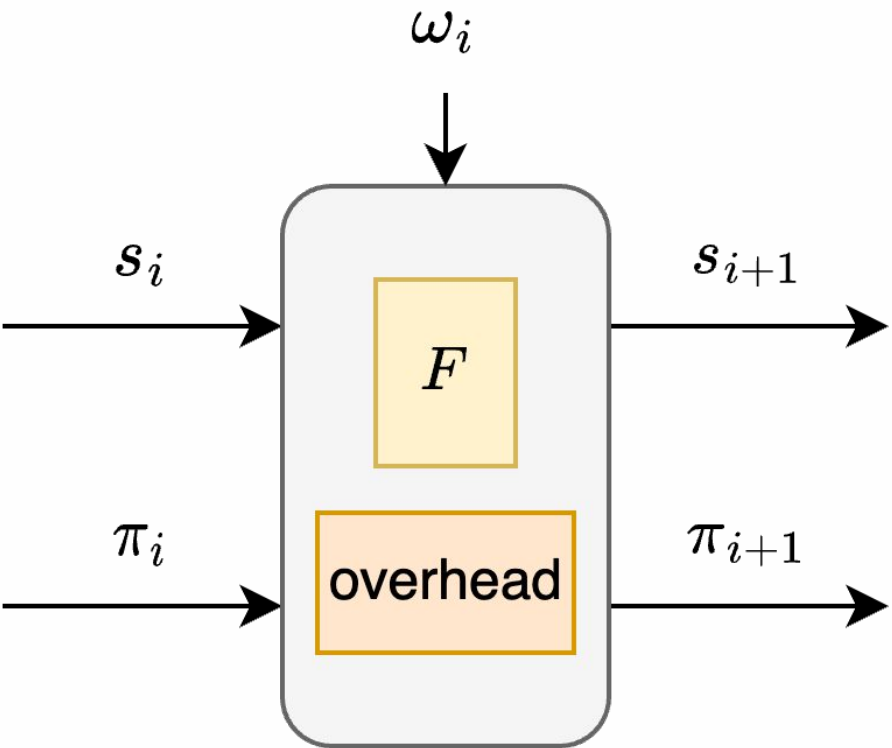


	Overhead	Examples
SNARK of SNARK	Read whole proof, run full verifier, defer nothing	Plonky2 , recursive STARK, Fractal , Groth16 + [BCTV14]
Accumulation (atomic)	Read whole proof, run partial verifier, defer hard verification	Halo/Halo2 , [BCMS20]



Recursion Overheads

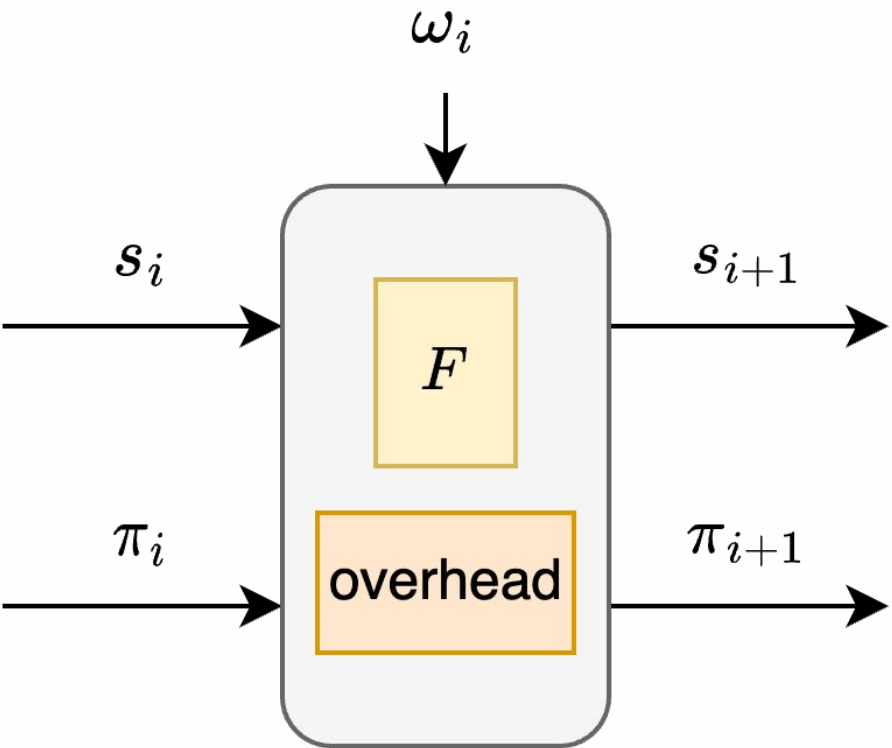
Recent works have drastically reduced the recursion overhead:



	Overhead	Examples
SNARK of SNARK	Read whole proof, run full verifier, defer nothing	Plonky2 , recursive STARK, Fractal , Groth16 + [BCTV14]
Accumulation (atomic)	Read whole proof, run partial verifier, defer hard verification	Halo/Halo2 , [BCMS20]
Accumulation (split)	Read partial proof, run partial verifier, defer hard verification and reading whole proof	[BCLMS21]

Recursion Overheads

Recent works have drastically reduced the recursion overhead:

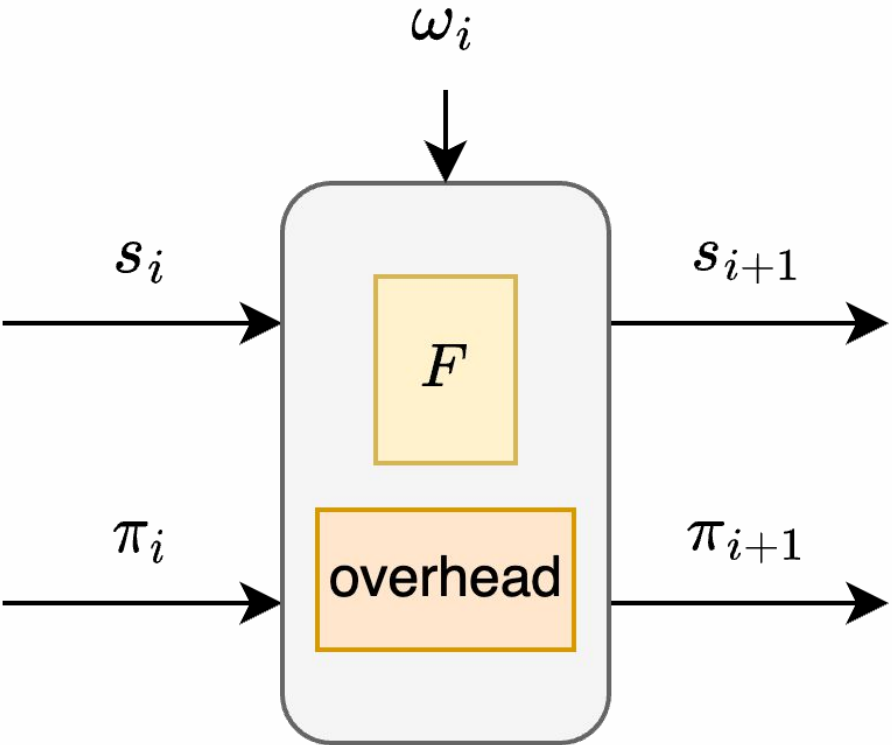


	Overhead	Examples
SNARK of SNARK	Read whole proof, run full verifier, defer nothing	Plonky2 , recursive STARK, Fractal , Groth16 + [BCTV14]
Accumulation (atomic)	Read whole proof, run partial verifier, defer hard verification	Halo/Halo2 , [BCMS20]
Accumulation (split)	Read partial proof, run partial verifier, defer hard verification and reading whole proof	[BCLMS21]
Folding	Read an <u>unproven</u> instance, compress it into a running instance, defer proving	Nova



Recursion Overheads

Recent works have drastically reduced the recursion overhead:



	Overhead	Examples
SNARK of SNARK	Read whole proof, run full verifier, defer nothing	Plonky2 , recursive STARK, Fractal , Groth16 + [BCTV14]
Accumulation (atomic)	Read whole proof, run partial verifier, defer hard verification	Halo/Halo2 , [BCMS20]
Accumulation (split)	Read partial proof, run partial verifier, defer hard verification and reading whole proof	[BCLMS21]
Folding	Read an <u>unproven</u> instance, compress it into a running instance, defer proving	Nova

Unlikely we can do better than deferring the proving. However, we can maybe have a more efficient circuit representation



Definition: PLONK Arithmetization

- PLONK traces are like Sudokus
 - fixed-size grid
 - fill the cells with “numbers”
 - set of rules to follow

q _L	q _R	q _O	q _M	q _C		a	b	c



Definition: PLONK Arithmetization

q _L	q _R	q _O	q _M	q _C	a	b	c
1	0	0	0	0	5		
1	0	0	0	0	10		
-1	0	0	1	0			
-1	0	0	1	0			
1	1	-1	0	0			
0	0	-1	1	0			

- PLONK traces are like Sudokus
 - fixed-size grid
 - fill the cells with “numbers”
 - set of rules to follow
- Some values are already filled in: these are the **selectors** and **public inputs**



Definition: PLONK Arithmetization

q _L	q _R	q _O	q _M	q _C		a	b	c
1	0	0	0	0		5		
1	0	0	0	0		10		
-1	0	0	1	0				
-1	0	0	1	0				
1	1	-1	0	0				
0	0	-1	1	0				

- PLONK traces are like Sudokus
 - fixed-size grid
 - fill the cells with “numbers”
 - set of rules to follow
- Some values are already filled in: these are the **selectors** and **public inputs**
- Rule #1: copy constraints



Definition: PLONK Arithmetization

q _L	q _R	q _O	q _M	q _C		a	b	c
1	0	0	0	0		5		
1	0	0	0	0		10		
-1	0	0	1	0				
-1	0	0	1	0				
1	1	-1	0	0				
0	0	-1	1	0				

- PLONK traces are like Sudokus
 - fixed-size grid
 - fill the cells with “numbers”
 - set of rules to follow
- Some values are already filled in: these are the **selectors** and **public inputs**
- Rule #1: copy constraints
- Rule #2: gate equation, applies at each row

$$(q_L)_i a_i + (q_R)_i b_i + (q_O)_i c_i + (q_M)_i a_i b_i + (q_C)_i = 0$$



Definition: PLONK Arithmetization

q_L	q_R	q_O	q_M	q_C			
1	0	0	0	0	5	10	
1	0	0	0	0			
-1	0	0	1	0			
-1	0	0	1	0			
1	1	-1	0	0			
0	0	-1	1	0			

$\mathbf{X} = (\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c)$

$\mathbf{W} = (\mathbf{w}_a, \mathbf{w}_b, \mathbf{w}_c)$

$$(\mathbf{q}_L)_i \mathbf{a}_i + (\mathbf{q}_R)_i \mathbf{b}_i + (\mathbf{q}_O)_i \mathbf{c}_i + (\mathbf{q}_M)_i \mathbf{a}_i \mathbf{b}_i + (\mathbf{q}_C)_i = 0$$

- PLONK traces are like Sudokus
 - fixed-size grid
 - fill the cells with “numbers”
 - set of rules to follow
- Some values are already filled in: these are the **selectors** and **public inputs**
- Rule #1: copy constraints
- Rule #2: gate equation, applies at each row
- Selectors and rules define a circuit
- Trace be divided into an instance and witness



Definition: Folding Schemes

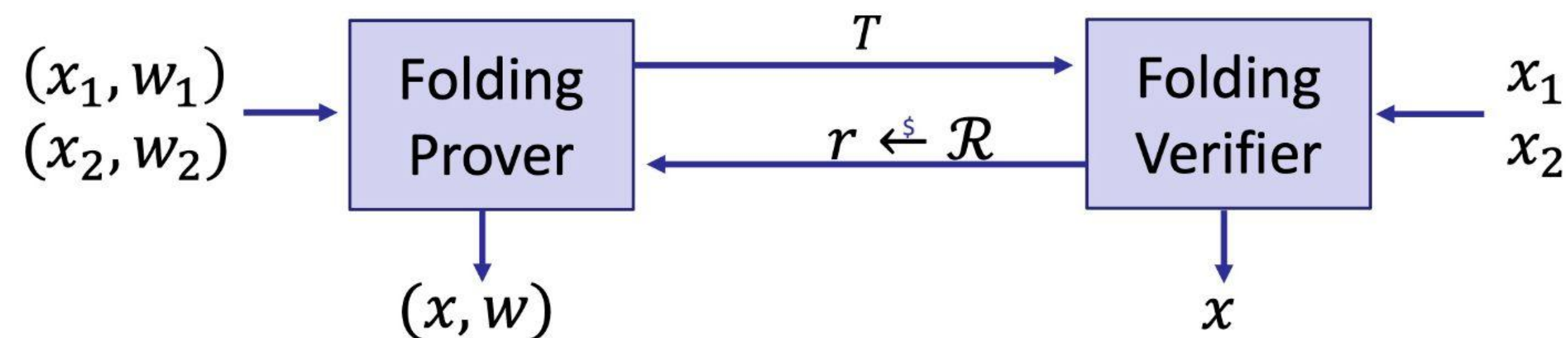


Definition: Folding Schemes

A folding scheme: compress two instances into one

Let $C: \mathbb{F}_p^n \times \mathbb{F}_p^m \rightarrow \mathbb{F}_p$ be a circuit

A folding scheme for C is a protocol between two parties:



Complete: if $C(x_1, w_1) = C(x_2, w_2) = 0$ then $C(x, w) = 0$

Knowledge sound: $\forall P^* \exists E$ s.t. $\forall x_1, x_2: P^*$ outputs valid w for $x \Rightarrow E$ outputs valid w_1, w_2

ZKP MOOC



A folding scheme for PLONK: Attempt #1

- We use the cryptographer's best friend: a random linear combination

$$\begin{array}{|c|c|c|} \hline \text{trace}' & & \\ \hline a'_1 & b'_1 & c'_1 \\ \hline a'_2 & b'_2 & c'_2 \\ \hline a'_3 & b'_3 & c'_3 \\ \hline a'_4 & b'_4 & c'_4 \\ \hline \end{array} + r \begin{array}{|c|c|c|} \hline \text{trace}'' & & \\ \hline a''_1 & b''_1 & c''_1 \\ \hline a''_2 & b''_2 & c''_2 \\ \hline a''_3 & b''_3 & c''_3 \\ \hline a''_4 & b''_4 & c''_4 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline \text{final_trace} & & \\ \hline a_1 & b_1 & c_1 \\ \hline a_2 & b_2 & c_2 \\ \hline a_3 & b_3 & c_3 \\ \hline a_4 & b_4 & c_4 \\ \hline \end{array}$$

- Copy constraints are preserved
- The gate equation is not linear, we are going to run into some trouble



A folding scheme for PLONK: Attempt #2

- Apply Nova's insights:
 - the gate equation can be relaxed
 - the Prover works over the trace while the Verifier works with commitments



A folding scheme for PLONK: Attempt #2

- Apply Nova's insights:
 - the gate equation can be relaxed
 - the Prover works over the trace while the Verifier works with commitments

- Relaxed PLONK arithmetization:

Witness: PLONK witness $\mathbf{W} = (\mathbf{w}_a, \mathbf{w}_b, \mathbf{w}_c)$ and an error vector \mathbf{e}

Instance: public inputs $\mathbf{X} = (\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c)$, a scalar u and commitments



Relaxed gate equation:

$$u [(\mathbf{q}_L)_i \mathbf{a}_i + (\mathbf{q}_R)_i \mathbf{b}_i + (\mathbf{q}_O)_i \mathbf{c}_i] + (\mathbf{q}_M)_i \mathbf{a}_i \mathbf{b}_i + u^2 (\mathbf{q}_C)_i + \mathbf{e}_i$$



A folding scheme for PLONK: Attempt #2

- Folding:

$$\text{final_trace} = \text{trace}' + r * \text{trace}''$$

$$u = u' + ru''$$

$$\mathbf{e} = ?$$



A folding scheme for PLONK: Attempt #2

- Folding:

$$\text{final_trace} = \text{trace}' + r * \text{trace}''$$

$$u = u' + ru''$$

$$\mathbf{e} = ?$$

- Plugging in the new values to the relaxed gate equation yields:

$$\text{Gate}(\text{final_trace}) = \text{Gate}(\text{trace}') + r\mathbf{t} + r^2 * \text{Gate}(\text{trace}'')$$



A folding scheme for PLONK: Attempt #2

- Folding:

$$\text{final_trace} = \text{trace}' + r * \text{trace}''$$

$$u = u' + ru''$$

$$\mathbf{e} = ?$$

- Plugging in the new values to the relaxed gate equation yields:

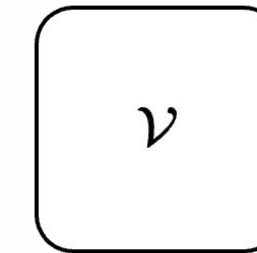
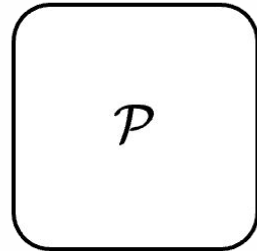
$$\text{Gate}(\text{final_trace}) = \text{Gate}(\text{trace}') + r\mathbf{t} + r^2 * \text{Gate}(\text{trace}'')$$

- We can get rid of \mathbf{t} by defining the error term and its folding as:

$$\mathbf{e} = \mathbf{e}' - r\mathbf{t} + r^2\mathbf{e}''$$



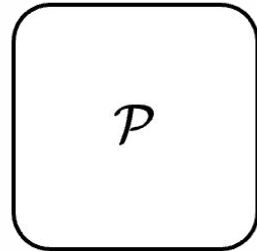
A folding scheme for PLONK: Sangria



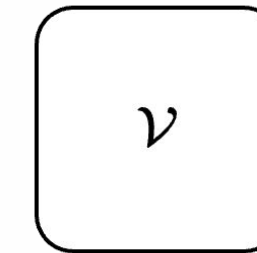
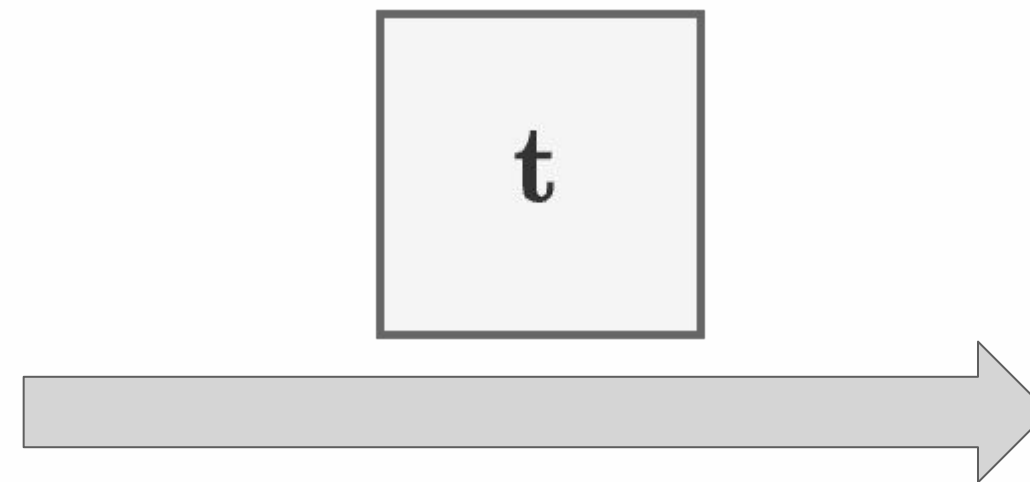
- Compute \mathbf{t} , as prescribed by distributing the gate equation
- Commit to \mathbf{t}



A folding scheme for PLONK: Sangria



- Compute \mathbf{t} , as prescribed by distributing the gate equation
- Commit to \mathbf{t}

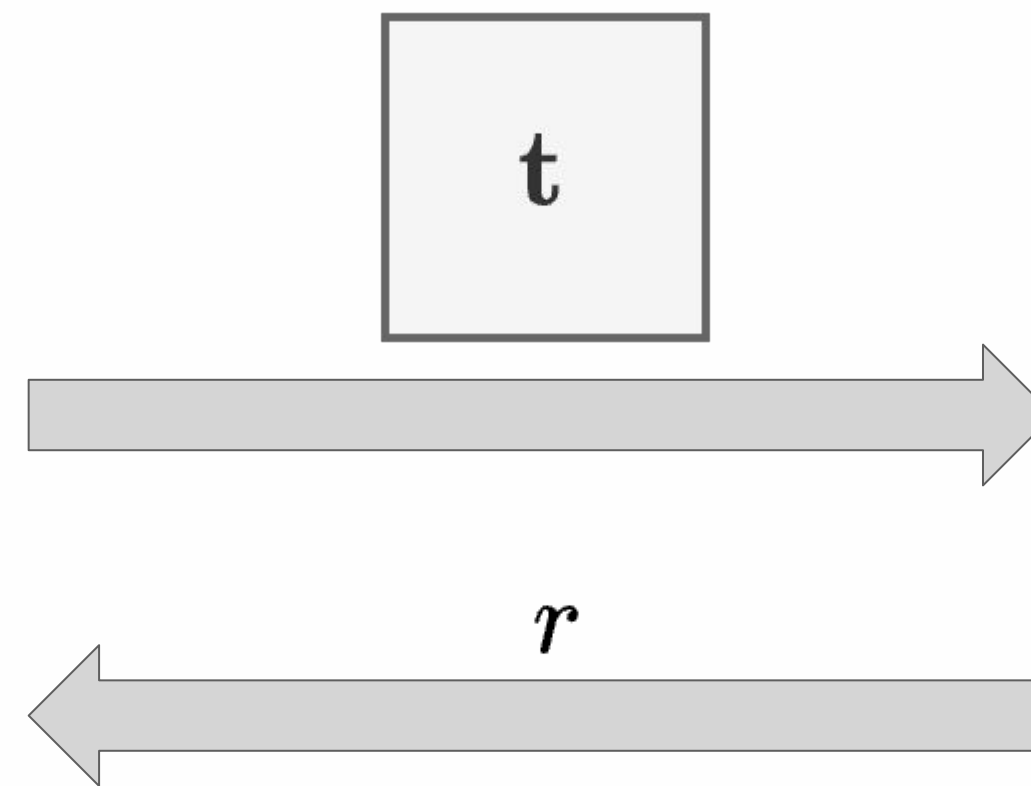


A folding scheme for PLONK: Sangria

\mathcal{P}

- Compute \mathbf{t} , as prescribed by distributing the gate equation
- Commit to \mathbf{t}

ν



- Sample a random value r



A folding scheme for PLONK: Sangria

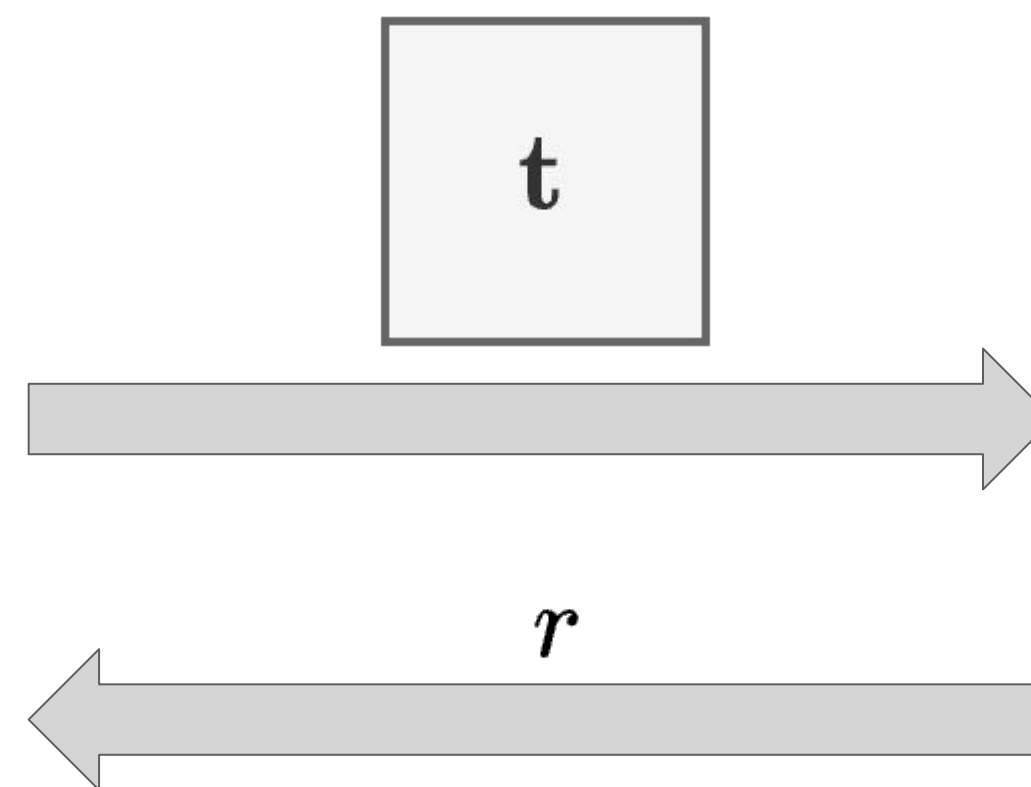
\mathcal{P}

- Compute \mathbf{t} , as prescribed by distributing the gate equation
- Commit to \mathbf{t}

- Compute linear combinations of the trace values

$$\text{final_trace} = \text{trace}' + r * \text{trace}''$$

$$\mathbf{e} = \mathbf{e}' - r\mathbf{t} + r^2\mathbf{e}''$$



\mathcal{V}

- Sample a random value r
- Compute linear combinations of the public inputs and commitments:

$$\begin{array}{lclcl} \boxed{\mathbf{w}_a} & = & \boxed{\mathbf{w}'_a} & + & r \boxed{\mathbf{w}''_a} \\ \boxed{\mathbf{w}_b} & = & \boxed{\mathbf{w}'_b} & + & r \boxed{\mathbf{w}''_b} \\ \boxed{\mathbf{w}_c} & = & \boxed{\mathbf{w}'_c} & + & r \boxed{\mathbf{w}''_c} \\ \boxed{\mathbf{e}} & = & \boxed{\mathbf{e}'} & - & r \boxed{\mathbf{t}} + r^2 \boxed{\mathbf{e}''} \end{array}$$



The verifier \mathcal{V} is given the verifier key \mathbf{vk} and two committed relaxed PLONK instances, $(\mathbf{X}', u', \overline{W'_a}, \overline{W'_b}, \overline{W'_c}, \overline{E'})$ and $(\mathbf{X}'', u'', \overline{W''_a}, \overline{W''_b}, \overline{W''_c}, \overline{E''})$. The prover \mathcal{P} is given the prover key \mathbf{pk} and both instances with their corresponding witnesses $(\mathbf{W}', \mathbf{e}', r'_a, r'_b, r'_c, r'_e)$ and $(\mathbf{W}'', \mathbf{e}'', r''_a, r''_b, r''_c, r''_e)$.

The Sangria folding scheme proceeds as follows:

1. \mathcal{P} samples r_t at random and sends $\overline{T} = \text{Com}(\text{pp}_E, \mathbf{t}; r_t)$ where \mathbf{t} is computed as:

$$\mathbf{t} := u''(\mathbf{q}_L \circ \mathbf{a}' + \mathbf{q}_R \circ \mathbf{b}' + \mathbf{q}_O \circ \mathbf{c}') + u'(\mathbf{q}_L \circ \mathbf{a}'' + \mathbf{q}_R \circ \mathbf{b}'' + \mathbf{q}_O \circ \mathbf{c}'') \quad (7)$$

$$+ \mathbf{q}_M \circ (\mathbf{a}' \circ \mathbf{b}'' + \mathbf{a}'' \circ \mathbf{b}') \quad (8)$$

$$+ 2ru'u''\mathbf{q}_C \quad (9)$$

2. \mathcal{V} samples the challenge r at random.
3. \mathcal{P} and \mathcal{V} output the folded instance $(\mathbf{X}, u, \overline{W_a}, \overline{W_b}, \overline{W_c}, \overline{E})$ where:

$$\mathbf{X} \leftarrow \mathbf{X}' + r\mathbf{X}'' \quad (10)$$

$$u \leftarrow u' + ru'' \quad (11)$$

$$\overline{W_a} \leftarrow \overline{W'_a} + r\overline{W''_a} \quad (12)$$

$$\overline{W_b} \leftarrow \overline{W'_b} + r\overline{W''_b} \quad (13)$$

$$\overline{W_c} \leftarrow \overline{W'_c} + r\overline{W''_c} \quad (14)$$

$$\overline{E} \leftarrow \overline{E'} - r\overline{T} + r^2\overline{E''} \quad (15)$$

4. \mathcal{P} outputs the folded witness $(\mathbf{W}, \mathbf{e}, r_a, r_b, r_c, r_e)$ where:

$$\mathbf{W} \leftarrow \mathbf{W}' + r\mathbf{W}'' \quad (16)$$

$$r_a \leftarrow r'_a + r \cdot r''_a \quad (17)$$

$$r_b \leftarrow r'_b + r \cdot r''_b \quad (18)$$

$$r_c \leftarrow r'_c + r \cdot r''_c \quad (19)$$

$$\mathbf{e} \leftarrow \mathbf{e}' - r\mathbf{t} + r^2\mathbf{e}'' \quad (20)$$

$$r_e \leftarrow r'_e - r \cdot r_t + r^2 \cdot r''_e \quad (21)$$



Performance

- The verifier's work is dominated by the additions of commitments

$$\begin{array}{rclcl} \boxed{w_a} & = & \boxed{w'_a} & + & r \boxed{w''_a} \\ \boxed{w_b} & = & \boxed{w'_b} & + & r \boxed{w''_b} \\ \boxed{w_c} & = & \boxed{w'_c} & + & r \boxed{w''_c} \\ \boxed{e} & = & \boxed{e'} & - & r \boxed{t} + r^2 \boxed{e''} \end{array}$$

- In practice we will be folding a standard PLONK instance into a relaxed PLONK instance:
 - the standard instance will not have an error term
 - the verifier will only have to perform **4 point additions**



Wide Circuits and Custom Gates (TurboPLONK)

- We can deal with **wider circuits** (circuits with larger fan-in fan-out at each gate) by committing to each extra column
 - Cost: The verifier will have to perform **1 extra point addition per trace column**



Wide Circuits and Custom Gates (TurboPLONK)

- We can deal with **wider circuits** (circuits with larger fan-in fan-out at each gate) by committing to each extra column
 - Cost: The verifier will have to perform **1 extra point addition per trace column**
- Dealing with **higher degree gates**:
 - A gate of degree d will be relaxed with powers of u up to u^d
 - the error's folding equation will be of the form

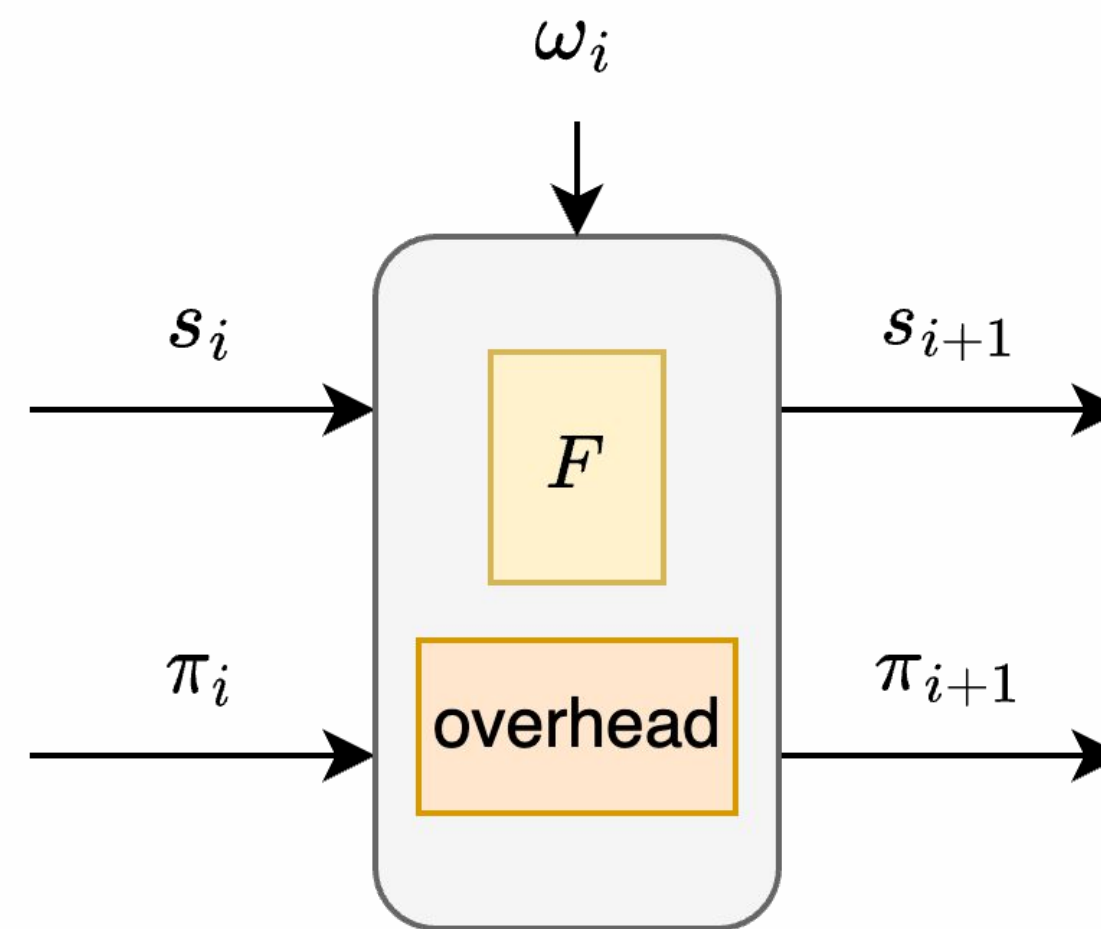
$$\mathbf{e} = \mathbf{e}' - r\mathbf{t}_1 - r^2\mathbf{t}_2 \cdots - r^{d-1}\mathbf{t}_{d-1} + r^d\mathbf{e}''$$

- Cost: each additional degree requires the Prover's message to include 1 new commitment. The Verifier will perform **1 extra point addition per degree**



New Design Space

Back to our original question, how small can we make the recursion overhead?



- Nova's circuit has approximately 20k constraints of which 12k are for point addition

Open question: can the benefit of wider circuits and custom gates outweigh the cost of the extra point additions that the folding verifier (recursion overhead) will perform?



Upcoming Work & Comments

- Sangria is a folding scheme for PLONK and piggybacks on Nova's results to achieve IVC. We can also piggyback on SuperNova!
- Implementation in progress for standard PLONK
- Folding lookup-enabled traces (ultraPLONK)
- (Hopefully) some circuit wizardry very soon, get in touch!





hello@geometry.xyz

@__geometry__



LONDON ST. HELIER SINGAPORE BELGRADE TEL AVIV BOSTON