

November 18, 2024

Succinct Arguments over Towers of Binary Fields

↑ Diamond & Posen 2023

Outline:

- ① Why SNARKs matter to industry - folk
- ② Overview of Modern SNARK Construction
- ③ Mathematical Preliminaries
 - Ⓐ Multilinear Polynomials
 - Ⓑ Towers of Binary Fields
- ④ The Binius IOP (Bird's Eye)
 - Ⓐ Zerochek
 - Ⓑ Virtual Polynomials
 - Ⓒ Add Gadget Example
- ⑤ Small Field Multilinear Poly Commit Scheme (Tease Only)

Goal: To introduce a technical audience to a SOTA SNARK proof system and shed light on what "cryptography" and "algorithms" industry cares for

Don't be shy to interrupt me with questions

⑤ Why SNARKs matter to industry-folks (Prelude)

Suppose you have

- Program P
- Input I

and you want to know Output $O := P(I)$

you generally have two options

* Run $P(I)$ yourself

Pro: Confident about the answer, 100% accuracy

Con: Expensive

100% cost

* Ask someone else to run for you

Pro: Free

0% cost

Con: No guarantee on output integrity

0% accuracy

SNARKs give us an almost perfect best of both worlds

* Ask someone else to run for you AND give you a SNARK proof

Pro: Close to free

< 0.1% cost

Pro: Great guarantee on output integrity

> 99.99...9% accuracy

SNARKs enable verifiable computation, particularly great when you have asymmetric computational power

Main Use Case Today

- Blockchains → Scaling Throughput

Ethereum is a very slow computer; Natively can validate 15 transactions per second
BUT... it can validate a SNARK proof of 10000 transactions very quickly

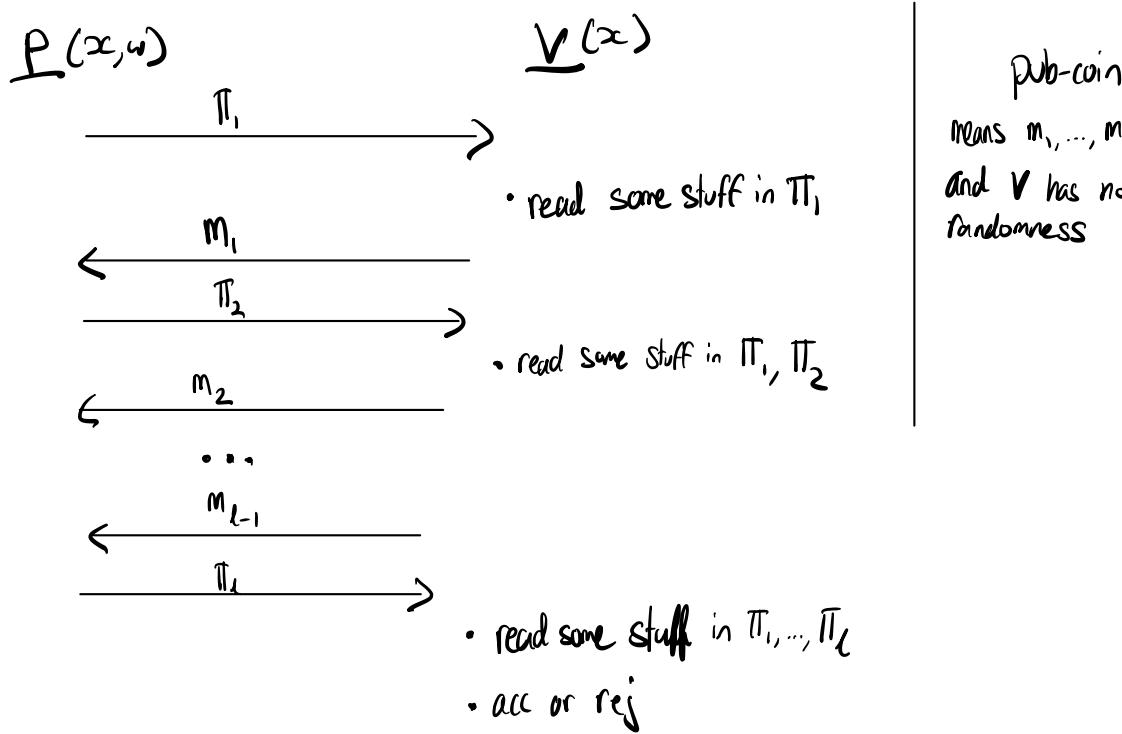
(Many other uses, but not our focus)

① Overview of Modern SNARK Construction

Key Tools:

- "public-coin" Interactive Oracle Proof (IOP)
- Commitment Scheme
- Hash function

What is an IOP?



Notes: Π_1, \dots, Π_l are oracles

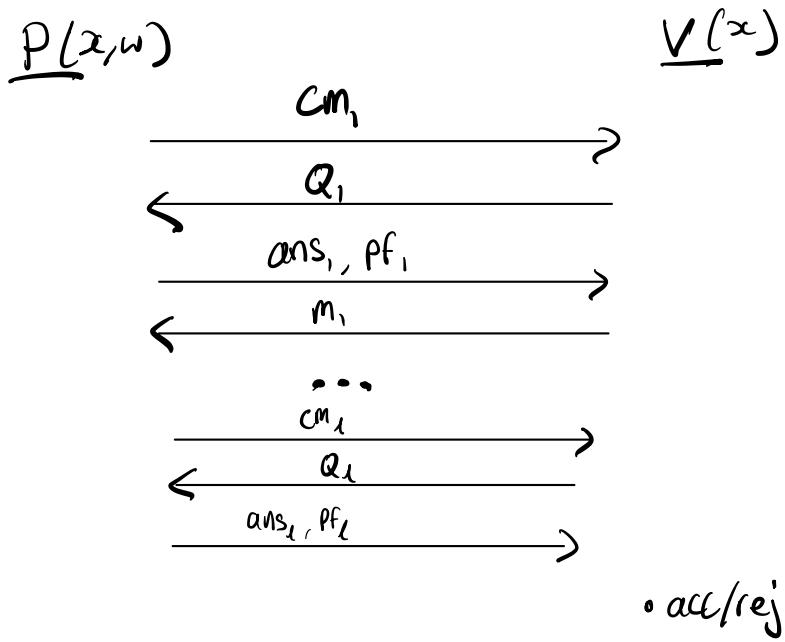
Much like PCPs, we can replace Π_i with an appropriate commitment c_i and answer verifier queries with opening proofs

Doing so, we compile an IOP into a succinct interactive proof
↑
if V isn't querying everywhere

IOP is good if:

- complete
- (knowledge)-sound
- succinct
- few rounds
- good prover efficiency
- good verifier efficiency

IOP + Commitment Scheme \Rightarrow IP



Note: different types of IOPs

- Simplest to think of Π_1, \dots, Π_l as vectors (so use VCS)
- In practice, we often want to use Π_1, \dots, Π_l as polynomials (so use PCS)

Back to Modern SNARK construction...

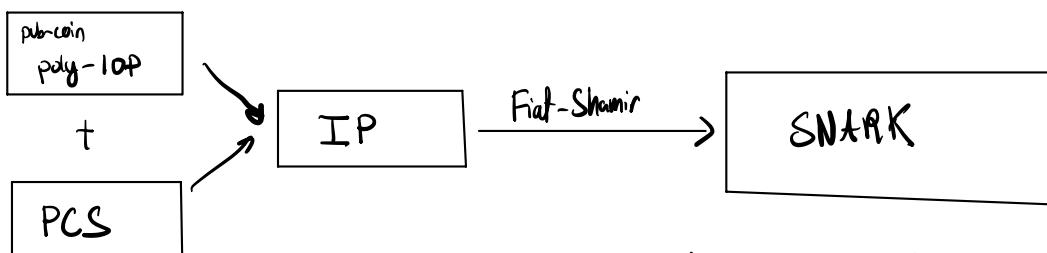
- Pick some NP-complete relation R (usually RICS or "Plonkish" relation)
- Construct a Poly-IOP for R (oracles are polynomials)
- Compile to IP with appropriate PCS

Front-End:

- Reduce arbitrary computation to an R -instance

[BCS16 transformation]

Back-End:



IP good + PCS good \Rightarrow IP good [CDGFS23]

2

Mathematical Preliminaries:

A Multilinear Polynomials

Defn: Multilinear polynomial is linear in each variable

$$\text{ex} \quad f(x_0, x_1, x_2) = x_0 + x_1 x_2 + x_0 x_1 + x_0 x_1 x_2 \quad \checkmark$$

$$\text{ex} \quad f(x_0, x_1, x_2) = x_0 + x_0^2 x_1 + x_1 x_2 \quad \times \quad (\deg_{x_0} f = 2 > 1)$$

Fact: A multilinear polynomial on n vars can be uniquely specified by its evaluations on the n -dimensional boolean hypercube

Given $f: \{0,1\}^n \rightarrow \mathbb{F}_3$ $\exists! \tilde{f} \in \mathbb{F}[x_0, \dots, x_{n-1}]^{\leq 1}$ st

$$\forall x \in \{0,1\}^n, f(x) = \tilde{f}(x)$$

Defn: We call \tilde{f} the multilinear extension of f

$$\text{ex} \quad \text{Sps } f(0,0) = a_{0,0}, f(1,0) = a_{1,0}, f(0,1) = a_{0,1}, f(1,1) = a_{1,1}$$

$$\begin{aligned} \text{Then } \tilde{f}(x_0, x_1) &= a_{0,0} L_{0,0}(x_0, x_1) \\ &\quad + a_{1,0} L_{1,0}(x_0, x_1) \\ &\quad + a_{0,1} L_{0,1}(x_0, x_1) \\ &\quad + a_{1,1} L_{1,1}(x_0, x_1) \end{aligned}$$

$\forall r_0, r_1 \in \mathbb{F}^2$

$$\text{evaluation: } \tilde{f}(r_0, r_1) = \sum_{b_0, b_1 \in \{0,1\}} a_{b_0, b_1} L_{b_0, b_1}(r_0, r_1)$$

$\uparrow \quad \uparrow$

Lagrange coeff. Lagrange Basis Polynomial.

$$L_{b_0, b_1}(x_0, x_1) = \begin{cases} 1, & \text{if } (x_0, x_1) = (b_0, b_1) \\ 0, & \text{if } (x_0, x_1) \in \{0,1\}^2 \setminus \{(b_0, b_1)\} \end{cases}$$

$$L_{0,0}(x_0, x_1) = (1-x_0)(1-x_1)$$

$$L_{1,0}(x_0, x_1) = x_0(1-x_1)$$

$$L_{0,1}(x_0, x_1) = (1-x_0)x_1$$

$$L_{1,1}(x_0, x_1) = x_0 x_1$$

↑ 4 Lagrange Basis Polynomials.

Why low degree extension useful? Short answer: amplifies small differences btwn \vec{a} vectors
 Long answer: Stay tuned

(B) Towers of Binary Fields

[Credit: Ben Diamond slides]

- Recall all finite fields of same size are isomorphic
- BUT we still want to choose our specific fields carefully.
↳ affects efficiency!

We choose a sequence of fields $\mathbb{F}_2 \subset \mathbb{F}_4 \subset \mathbb{F}_{16} \subset \dots \subset \mathbb{F}_{2^{128}}$
 \parallel
 $T_0 \subset T_1 \subset T_2 \subset \dots \subset T_7$

where $T_i \cong \mathbb{F}_{2^{2^i}}$ and can be represented as a 2^i bit-string

Formally:

$$\left\{ \begin{array}{l} T_0 := \mathbb{F}_2 \\ T_1 := T_0[X_0] / (X_0^2 + X_0 + 1) \cong \mathbb{F}_4 \\ T_2 := T_1[X_1] / (X_1^2 + X_0 X_1 + 1) \cong \mathbb{F}_{16} \\ T_3 := T_2[X_2] / (X_2^2 + X_1 X_2 + 1) \cong \mathbb{F}_{256} \\ \dots \\ T_{l+1} := T_l[X_l] / (X_l^2 + X_{l-1} X_l + 1) \end{array} \right.$$

Wiedemann Tower [1987]

$X_l^2 + X_l X_{l-1} + 1$
always irreducible

Fan-Paar 1997: efficient multiplication and inversion in Wiedemann's Tower

- $O(n^{1.59})$ where $n = 2^l$ the bit-width (Karatsuba-like)

Why not use classical approach: $\mathbb{F}_{2^{128}} \cong \mathbb{F}_2[X] / (X^{128} + \dots)$

- Support multiple data types (data widths) in arithmeticization
- RS codes alphabets are intermediate size $\mathbb{F}_{2^{16}}$
- Intermediate fields need to be efficient

optimal in tower case
costs in classical case

• efficiently "packable": assemble 16 \mathbb{F}_2 values into 1 $\mathbb{F}_{2^{\ell}}$ value
 {
 • efficient "small" by "large" mult:
 $\mathbb{F}_{2^{\ell}} \times \mathbb{F}_{2^{128}}$ mult FASTER than $\mathbb{F}_{2^{128}} \times \mathbb{F}_{2^{128}}$
 • efficiently "unpackable"

$\mathbb{F}_{2^{\ell}}$ always a
16-dimensional
 \mathbb{F}_2 -vector space
but identifiable?

- In classical approach you see 128 bits for some $\mathbb{F}_{2^{128}}$ element

$$[10110010 \dots 1] \cdot [X^0, \dots, X^{127}] \rightarrow \text{some degree 127 poly}$$

is it in subfield $\mathbb{F}_{2^{64}}$? need to run some computation

- In tower approach, you see 128 bits for some $\mathbb{F}_{2^{128}}$ element

$$[101100100000 \dots 0] \cdot [1, X_0, X_1, X_0X_1, X_2, X_0X_2, X_1X_2, X_0X_1X_2, X_3, \dots, X_0X_1X_2X_3X_4X_5X_6]$$

gives some multilinear in 7 variables

- * The first 64 bits do not contribute to any monomial terms with X_6
 * The last 64 bits all contribute to monomials terms with X_6

Element in $\mathbb{F}_{2^{64}}$ \longleftrightarrow last 64 bits are all 0.

Similarly can pack 2 $\mathbb{F}_{2^{64}}$ values into 1 $\mathbb{F}_{2^{128}}$ value using first/last 64 bits

Generally T_{l+k} is a 2^k -dimensional T_l -vector space with basis elements

$$[1, X_l, X_{l+1}, X_lX_{l+1}, X_{l+2}, \dots, X_lX_{l+1} \dots X_{l+k-1}]$$

- So, small-by-large mult is 2^k small-by-smalls
 (good b/c recall same-size mult is super-linear, $O(n^{1.59})$)

$$\text{cost of } 128 \times 128 \text{ mult} = 128^{1.59} = 32^{1.59} \cdot 4^{1.59} > 32^{1.59} \cdot 4 = \text{cost of } 4 \text{ } 32 \times 32 \text{ mults} = \text{cost of } 128 \times 32 \text{ mult}$$

③ Bird's Eye Perspective of the Binus IOP

- Planckish Arithmetization is our chosen NP-complete relation

Instance:

15	251			
		42		
			26	
				7870

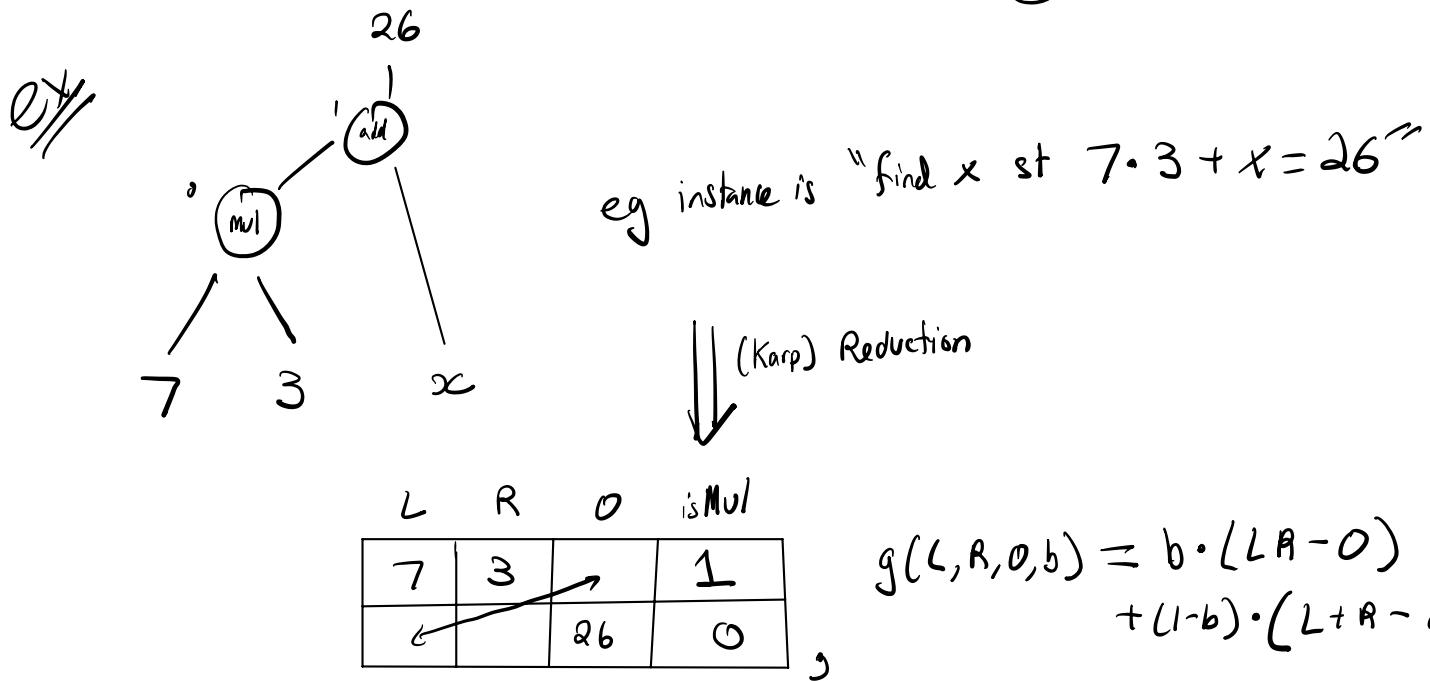
- a partially filled out 2D table
- some copy constraints indicates two cells must have same value
- a gate polynomial of arity #cols, g

Witness:

a fully filled out table such that

- agrees with partial filling
- copy constraints hold
- $\forall \vec{row}_i, g(\vec{row}_i) = 0 \quad // \text{gate constraints hold.}$

- Can reduce from Arithmetic Circuit-SAT (Circuit sat over some field, not necessarily \mathbb{F}_2)
"even more" NP-Hard $\textcircled{11}$



In practice: most front-ends turn arbitrary computations into some \mathbb{F}_p -arithmetic circuit (for large prime p) and then (kind of) apply this generic reduction

This sucks... every cell is treated as an \mathbb{F}_p -element

Usually most columns are 1s and 0s

Binus arithmeticization \rightarrow allows you to choose T_i for each column

Better for many reasons:

- i) Witness smaller in memory \rightarrow bandwidth is bottleneck for Hardware acceleration
- ii) Eventually we "commit to each column" \rightarrow (non-trivial, this work) it is more efficient to commit to n small values than n big values

A) Zerochek (for Gate Constraints)

Setting: I want to prove to you that the 10th fibonacci is 34

A	A'	A''
0	1	
		34

Instance:

Gate constraint:

$$g(A, A', A'') = A'' - A' - A$$

want for each row, $g(\text{row}) = 0$

"Copy" Constraints

$$\forall i: A[i+2] = A'[i+1] = A''[i]$$

(\Leftrightarrow diagonals are equal)

witness:

lab ... nib

	A	A'	A''
(0,0,0)	0	1	1
(1,0,0)	1	1	2
(0,1,0)	1	2	3
(1,1,0)	2	3	5
(0,0,1)	3	5	8
(1,0,1)	5	8	13
(0,1,1)	8	13	21
(1,1,1)	13	21	34

$$D := g(A, A', A'')$$

0	$\leftarrow g(0, 1, 1)$
0	$\leftarrow g(1, 1, 2)$
0	$\leftarrow g(1, 2, 3)$
0	
0	
0	
0	
0	$\leftarrow g(13, 21, 34)$

(the witness is also called the trace)

↑ gate evaluations

$$\vec{r} = (r_0, r_1, r_2) \quad \tilde{\lambda}(\vec{r}), \tilde{\lambda}'(\vec{r}), \tilde{\lambda}''(\vec{r})$$

$$\tilde{D}(\vec{r}) \leftarrow g(\tilde{\lambda}(\vec{r}), \tilde{\lambda}'(\vec{r}), \tilde{\lambda}''(\vec{r}))$$

Goal: Show that vector D is all zeroes

But that would require Verifier to know full witness (computation trace)

Idea: Pretend \tilde{D} is multilinear polynomial on 3-vars where the

$\vec{v} \in \{0, 1\}^3$, th entry $D[\vec{v}]$ is evaluation of $\tilde{D}(\vec{v})$

Then zero-test $\tilde{D}(x_0, x_1, x_2)$ @ random point $\vec{r} = (r_0, r_1, r_2) \in (\mathbb{F}_{2^{128}})^3$

$$\tilde{D}(\vec{r}) = 0 \Rightarrow \text{whp } D = \vec{0}$$

Why fast to evaluate $\tilde{D}(\vec{r})$?

- Poly commitment + opening proofs are succinct and efficiently verifiable

- Prover sends verifier commitments to all columns as multilinear polynomials.

↑ (not all, you will see)

Defn: Zerochek is an IOP for a claim that a specific multivariate polynomial evaluates to 0 over the hypercube

Here, take $h(x_0, x_1, x_2) := g(\tilde{A}(x_0, x_1, x_2), \tilde{A}'(x_0, x_1, x_2), \tilde{A}''(x_0, x_1, x_2))$

P shows V that $\forall \vec{v} \in \{0,1\}^3, h(\vec{v}) = 0$

Common mistake:

NOT proving $h(x_0, x_1, x_2) = 0$ (b/c h not necessarily multilinear
eg $g(A, A^2, A^3) = A^n - A^3 A$

However

$$\tilde{D}(x_0, x_1, x_2) = \sum_{\vec{v} \in \{0,1\}^3} L_{\vec{v}}(x) \cdot \underbrace{h(v)}_{\text{constant}}$$

is the multilinear poly whose hypercube evals, by design,
agree with h .

We check that $\tilde{D}(x_0, x_1, x_2) = 0$

I will stop unravelling Zerochek here.

On YouTube I have 2 uploaded videos about

- Zerochek and zerochek \rightarrow sumcheck reduction
- the sumcheck protocol

that are presented in the context of Binus. I encourage you to check those out

OK... so we have a way to prove gate constraints.

How about copy constraints?

Binus doesn't do copy constraints exactly... --

(B)

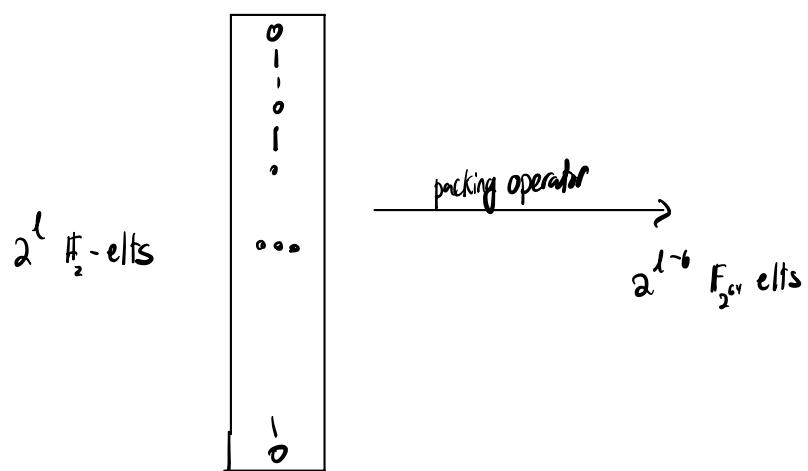
Virtual Multilinears:

Roughly

a virtual multilinear is a way to manipulate / "reshape" committed multilinear in some useful way. we highlight two

• Packing Construction

rows committed column A



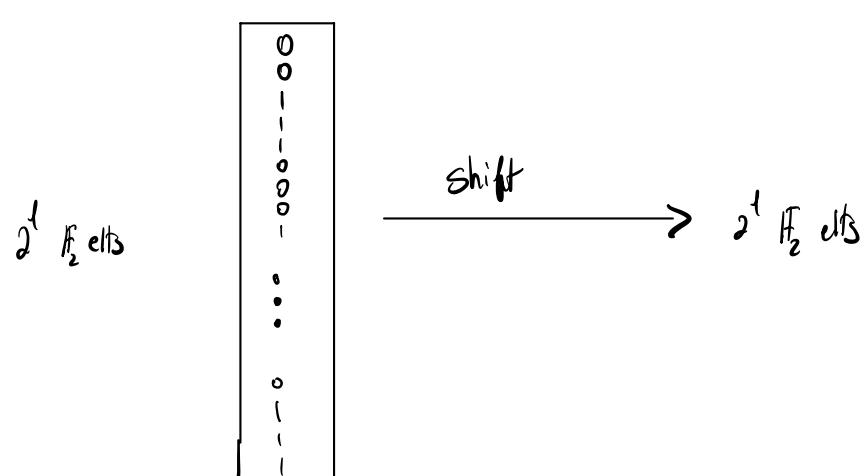
Pack₆(A)

Recall $F_{2^{64}}$ is a 64-dim \mathbb{F}_2 -vector space, with basis elts B_v for $v \in \{0,1\}^6$

$$\widehat{\text{pack}}_6(A)(x_0, \dots, x_{l-7}) = \sum_{v \in \{0,1\}^6} \tilde{A}(\vec{v}, x_0, x_1, \dots, x_{l-7}) \cdot B_v$$

• Shift Construction

rows committed column A



Virtual
Shifted Column

} Shifted by 3 positions
say.

(Can be circular or zero-fill, can be blockwise)

can
ignore

$$\text{Shift}_{b,0}(A)(x_0, \dots, x_{l-1}) = \sum_{v \in \{0,1\}^b} \hat{A}(\vec{v}, x_b, \dots, x_{l-1}) \cdot \tilde{s\text{-ind}}(\vec{v}, x_0, \dots, x_{b-1})$$

Shift lets us compare cols without committing to something new.

No need for copy-constraints in Fib example!

Commit to A and let A^1, A^2 be virtual shifts of A

Problem 1: Shift doesn't work? Fine modify like so

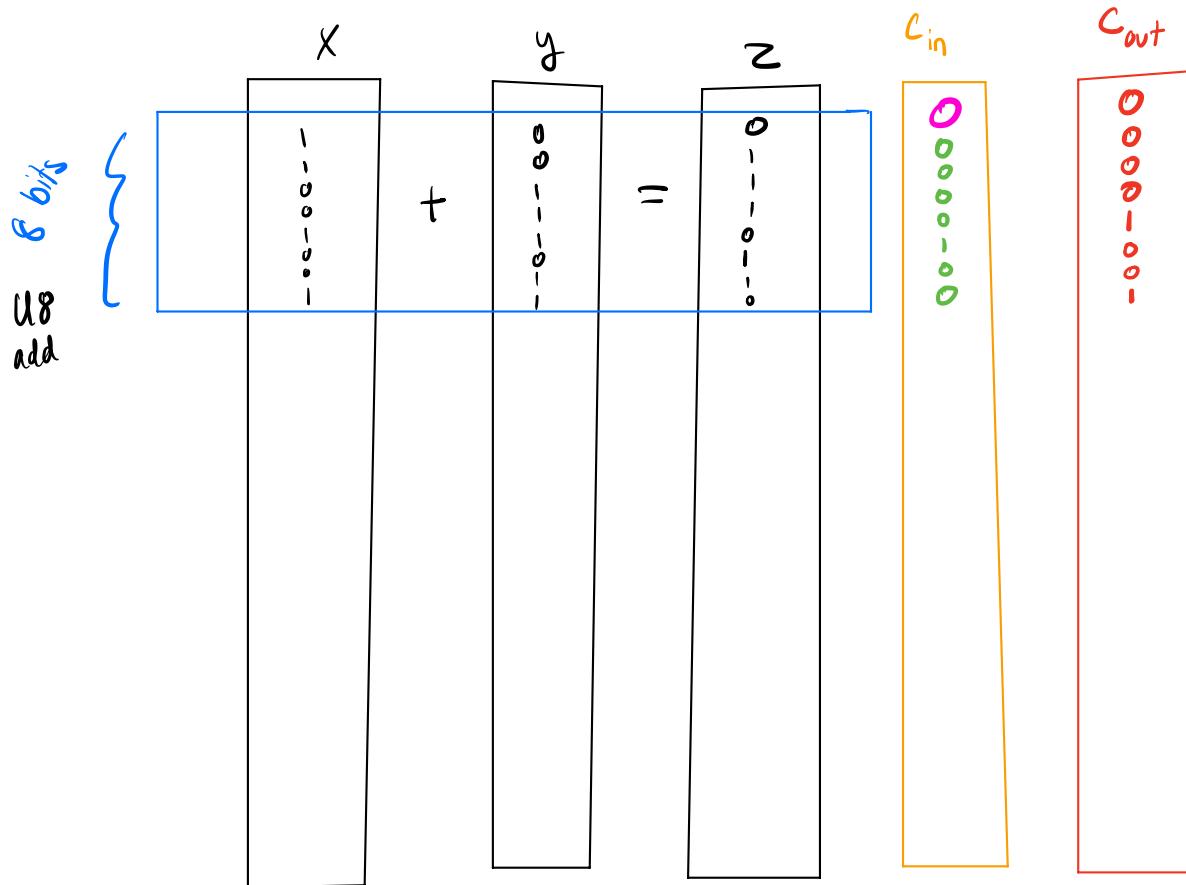
A	$\text{shift}_1(A)$	$\text{shift}_2(A)$	B
0	1	1	1
1	1	2	1
1	2	3	1
2	3	5	1
3	5	8	1
5	8	13	1
8	13	21	1
13	21	34	1
21	34	0	0
34	0	0	0

- highlighted values are instance
- Commit to A only

$$g(A, A^1, A^2, B) = B \cdot (A^2 - A^1 - A)$$

Problem 2: Subtraction is happening in binary field (just XOR of bits)
no good!

⑥ Add Gadget Example — Just a few \mathbb{F}_2 operations



- Commit C_{out}
- Declare $C_{in} = \text{shift-down of } C_{out} \text{ by 1 position, with zero-fill, block-size } 2^3 = 8$
- Constrain:

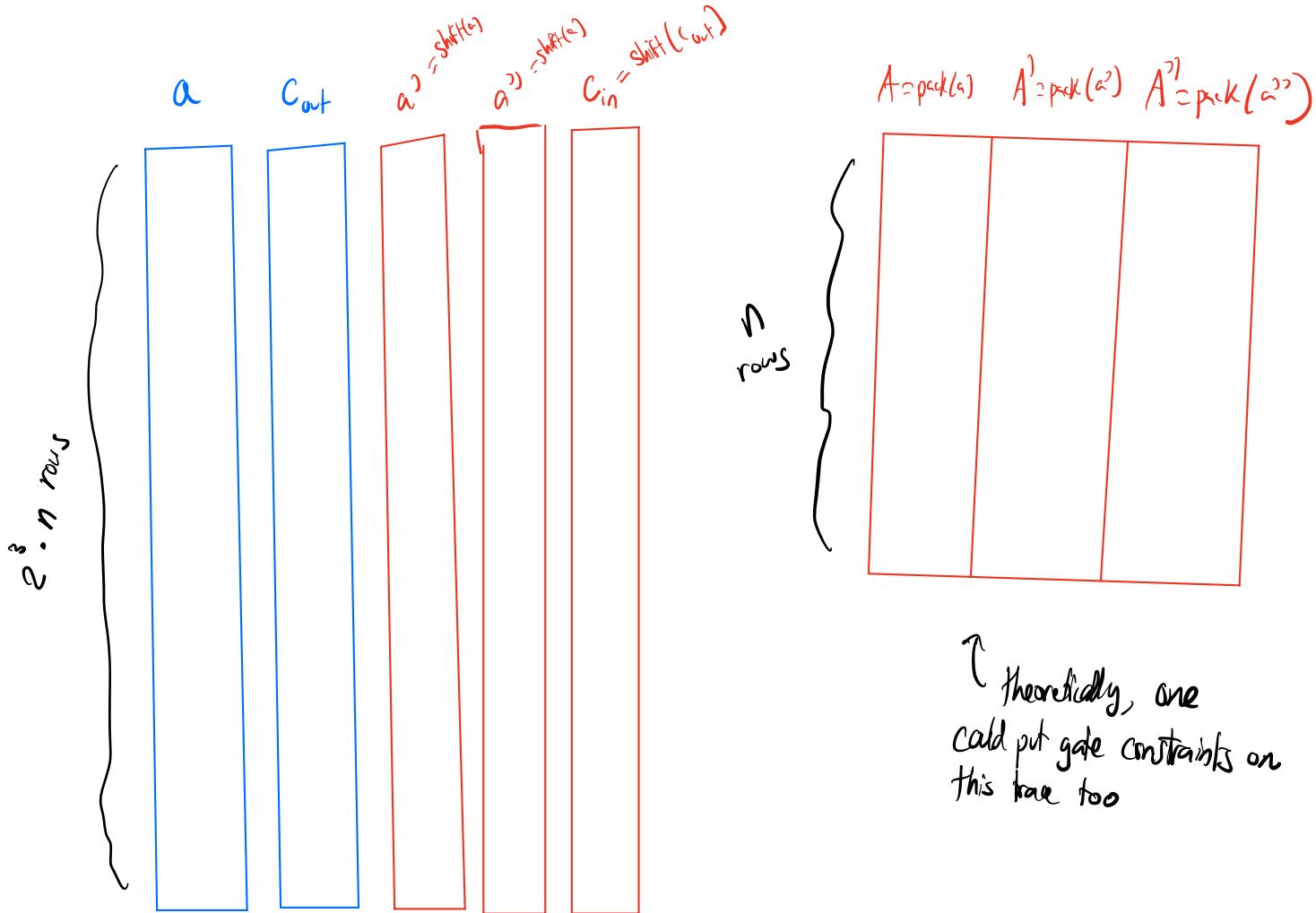
$$\bullet Z = X + Y + C_{in} \quad (\text{over } \mathbb{F}_2)$$

$$\bullet C_{out} = XY + XC_{in} + YC_{in} \quad (\text{over } \mathbb{F}_2) \quad \begin{matrix} C_{out}=1 \\ \text{if } \geq 2 \text{ of } \\ X, Y, C_{in}=1 \end{matrix}$$

(multiple constraints OK, batch into single constraint with random linear combination)

Just 2 constraints, shift, and extra column.

What does fibonacci look like now? (Ignore problem 1)



\nearrow
 theoretically, one
 could put gate constraints on
 this trace too

\nwarrow
gate constraints on this trace

$8n$ \mathbb{F}_2 elements

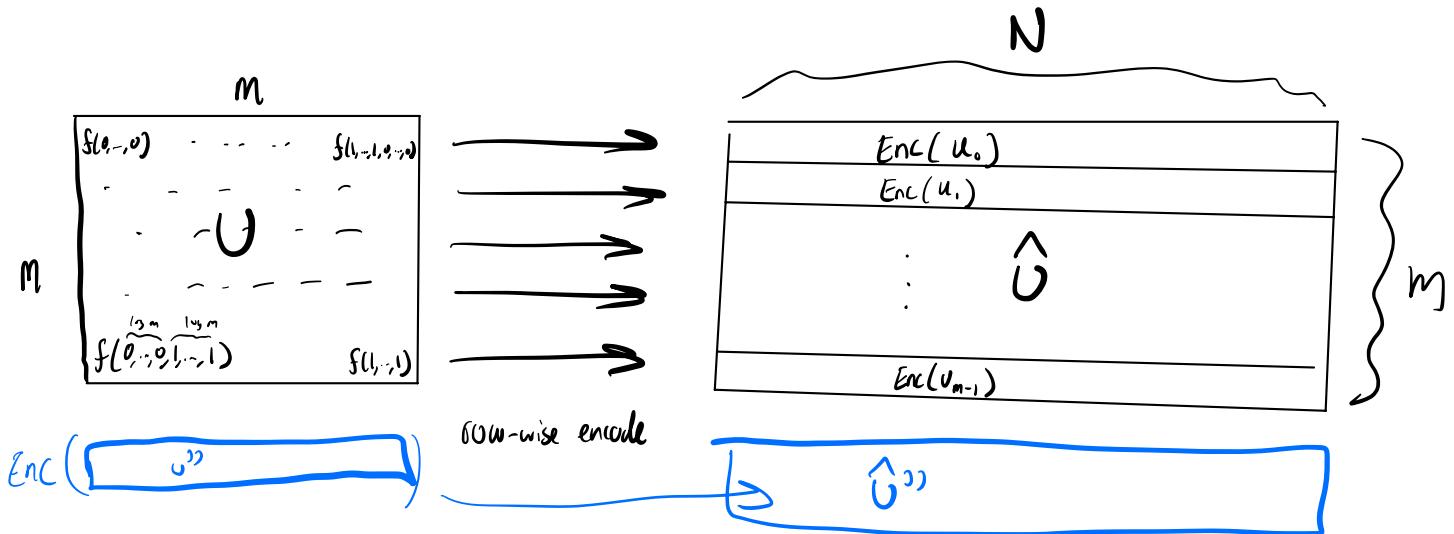
$n \quad \mathbb{F}_2^3 = \mathbb{F}_2^8 = \mathbb{F}_{256}$ elements

Commit to a and C_{out}

(4) Small-Field Multilinear Poly Commit Scheme

Let's recall the Brakedown PCS (Simplified)

- Let C be systematic linear code over \mathbb{F}_q with constant rate and relative distance
- Let $m = \sqrt{n}$ (for simplicity say m power of 2)
- Let $f(x_0, \dots, x_{l-1})$ be multilinear to commit where $l = \log n$



Evalu

To eval $f(r)$ for some $r \in \mathbb{F}^{\log n}$:

- Let $q_1, q_2 \in \mathbb{F}_q^m$ st $q_i = \bigotimes_{j=0}^{\log n-1} [1-r_j, r_j]$

$$q_2 = \bigotimes_{i=\log m}^{\log n-1} [-r_i, r_i]$$

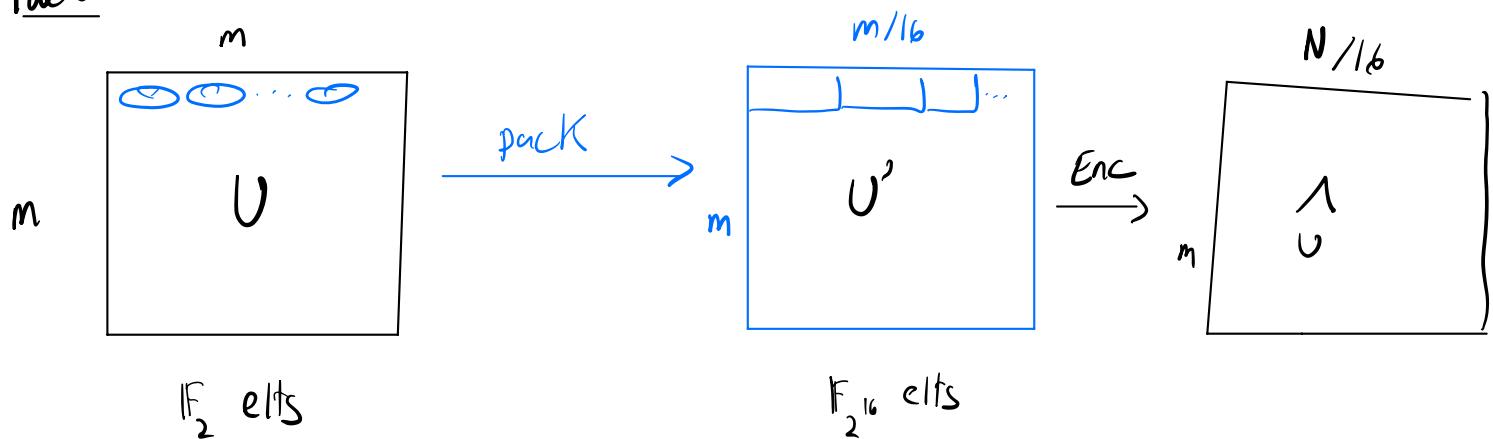
- P sends u^0 to V claimed to be q_1 -linear combination of V rows
- V checks consistency b/w $Enc(v^0)$ and q_1 -linear combination of \hat{v} (at $O(\lambda)$ columns only because code has good distance)
- V computes $f(r) = \langle u^0, q_2 \rangle$

If $f \in \mathbb{F}_{2^m}[x_0, \dots, x_{\ell-1}]^{\leq 1}$, then use normal Brakedown

what if $f \in \mathbb{F}_2[x_0, \dots, x_{\ell-1}]$?

- Do not want $16 \times$ blowup in size
- Maybe also wanted to constrain all values are in \mathbb{F}_2

Idea:



Commitment is \hat{U} (merkle root)

To eval @ $r = (r_0, \dots, r_{\ell-1}) \in (\mathbb{F}_{2^{128}})^{\ell}$

- P sends $u^0 = q_1\text{-linear combination of } U$ (unpacked)
- V checks consistency b/w $\text{Enc}(\text{pack}(u^0))$ and $q_1\text{-l.c. of } \hat{U}$ at $O(\lambda)$ column blocks
- V computes $f(r) = \langle u^0, q_2 \rangle$

Okay, I can pack \mathbb{F}_2 elements into $\mathbb{F}_{2^{16}}$ elements, easy

How on earth do I pack $\mathbb{F}_{2^{128}}$ elements into $\mathbb{F}_{2^{16}}$ elements ?!?

We don't actually use RS codes over the $\mathbb{F}_{2^{16}}$ alphabet,
our symbols are actually "tower algebra" elements 

Formally

$$A_{\iota, \kappa, \tau} := T_\tau [Y_0, \dots, Y_{\kappa-1}] / \left(Y_0^2 + X_{\iota-1} Y_0 + 1, Y_0^2 + Y_0 Y_1 + 1, \dots, Y_{\kappa-1}^2 + Y_{\kappa-2} Y_{\kappa-1} + 1 \right)$$

In Picture:

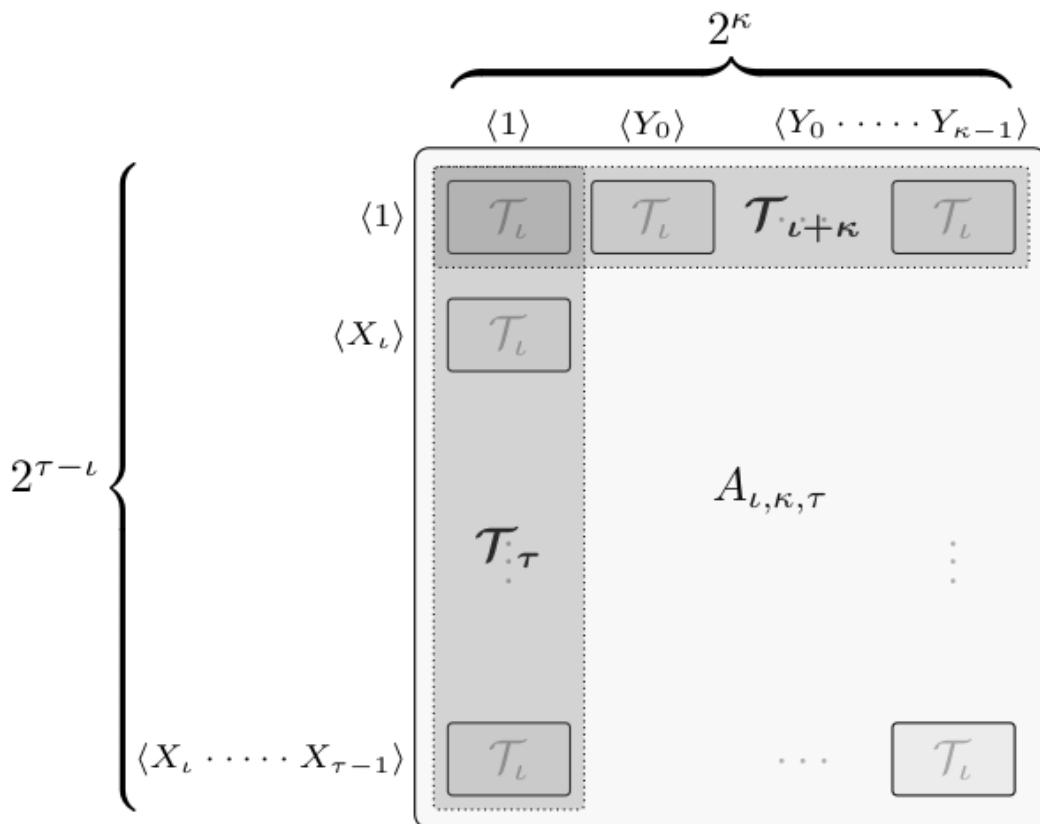


Figure 1: A depiction of our “tower algebra” data structure.

- So bits get packed into what is essentially 16-bit elements
- but big field elts get packed into a fully general $A_{\iota, \kappa, \tau}$ element

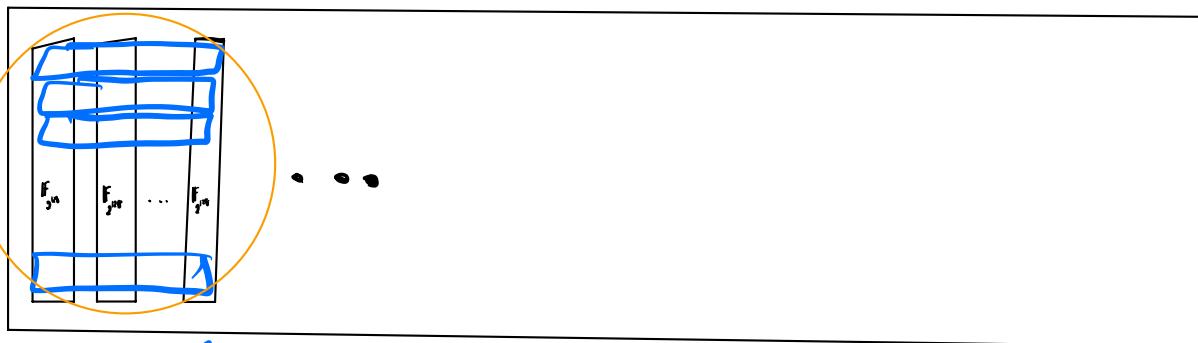
can do

$$\bullet T_\tau \times A_{\iota, \kappa, \tau}$$

and $\bullet A_{\iota, \kappa, \tau} \times T_{\iota+\kappa}$

So how do we check consistency?

$$v'' \in (\mathbb{F}_{2^{128}})^m$$



one tower algebra element, then encode these.

Compare $\text{Enc}(\text{embed}(v''))$ with q_1 -l.o.c. of \hat{G} rows at random cols. Fußk
 q_1 elts are $\mathbb{F}_{2^{128}} = T_L$, \hat{G} symbols are $A_{L,K,L}$ elts
so can left-multiply and everything type checks

of note: Need a new proximity gap result for optimal efficiency

TLDR the binding holds b/c of a statement that says

if one row of \hat{G} is δ -far from code

then random l.o.c. of \hat{G} rows is δ -far from code whp

(Proximity Gap survey would need a whole new talk)

want δ high so we can do fewer column-consistency checks to achieve λ bits of security

\Rightarrow smaller proof sizes