

LatticeFold & its Applications to Succinct Proof Systems

Dan Boneh

Binyi Chen



Stanford University

(zk)SNARKs

E.g., $\text{SHA-3}(\mathbf{w}) = \mathbf{x}$

(zk)SNARK = A *succinct* ZK proof showing that $\exists \mathbf{w}$ s.t. $C(\mathbf{x}, \mathbf{w}) = 0$

$\mathbf{S}_{(C)} \rightarrow (pp_C, vp_C)$



Properties:

- **Completeness:** honest P can compute valid π
- **Knowledge soundness:** malicious P^* knows valid \mathbf{w} if it can generate valid π
- **Zero knowledge:** π hide the witness \mathbf{w}

Key requirements for π : **Short** (i.e. $|\pi| \ll |\mathbf{w}|$) + **Fast to verify** (e.g. $O(\log|F|)$ time)

Applications: Blockchain, Verifiable zkML/FHE, Fighting disinformation & more

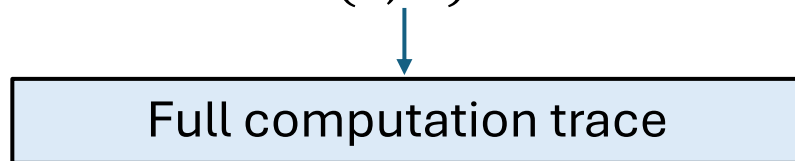
[Xie+22, NT16, DB22, KHSS22, BBBF18, XCBFCK22.....]

Challenges: Proving *expensive* statements (e.g., ML tasks) *efficiently*

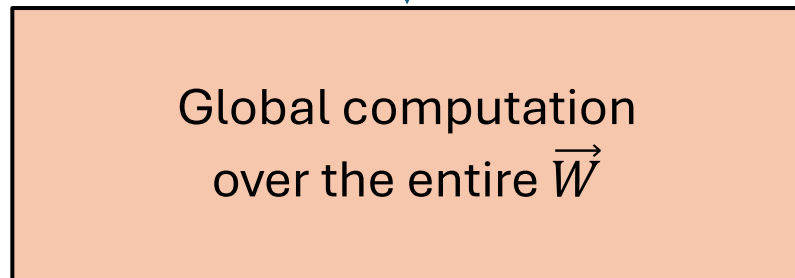
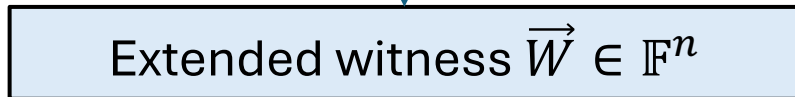
Monolithic SNARKs [\[Bitansky-Canetti-Chiesa-Tromer12...\]](#)

E.g., a block of 10K txs is valid w.r.t. ledger state update

NP statement (x, w) for a relation R_C



Algebraic transform



FFTs, MSMs, etc...

Proof π

Pre-quantum Schemes:

- Groth16, Plonk [\[GWC19\]](#), Marlin[\[CHMMVW20\]](#), Bulletproof[\[BBBPWM18\]](#)
- HyperPlonk[\[CBBZ22\]](#), Spartan[\[Setty19\]](#), etc...

Post-quantum Schemes:

- STARK[\[BBHR18\]](#), Brakedown[\[GLSTW21\]](#), Ligerio[\[AHIV17\]](#), Basefold[\[ZCF23\]](#) ...
- Lattice Bulletproofs[\[BLNS20,ACK21\]](#), LaBRADOR[\[BS22\]](#) ...

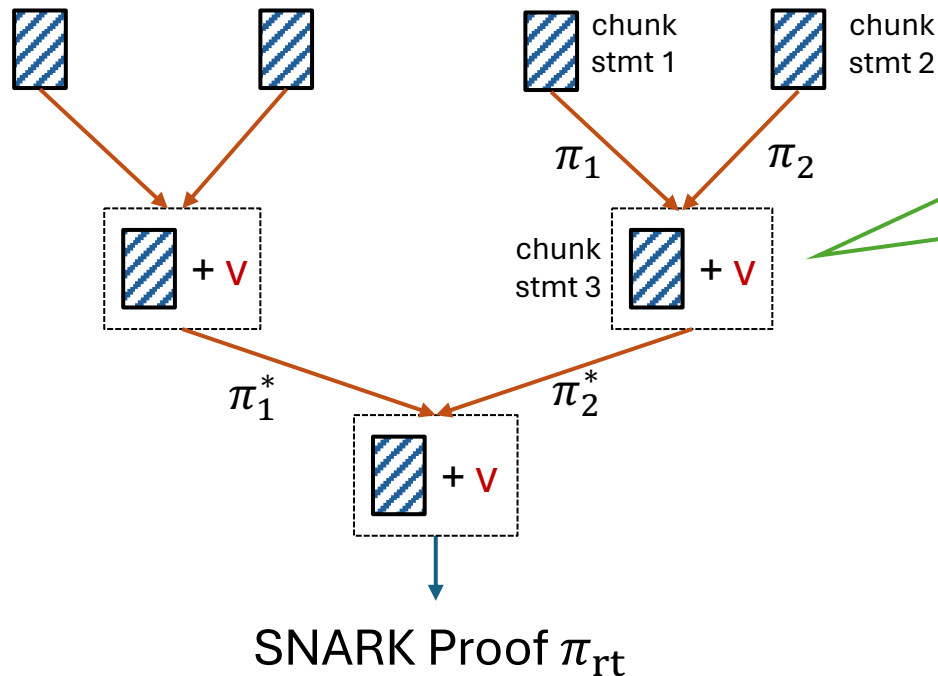
Challenges for proving **expensive** computation:

- Expensive global computation
- Large prover memory
- Harder parallelization + less streaming-friendly

Piecemeal SNARKs [Valiant08, BCTV14, BCCT12]

NP statement (x, w) for a relation R_C
e.g., a block of 10K txs is valid w.r.t. the ledger state

↓ split




Ideas:

e.g. Mangrove[NDCTB24], [Sou23]

- Split the statement into multiple small chunks
- Prove chunk statements using **SNARK Recursion** [Bitansky-Canetti-Chiesa-Tromer12]

SNARK Circuit:

- (1) chunk stmt 3  is correct
- (2) π_1, π_2 **verify** correctly

Pros:

- Minimal memory overhead
- Streaming/parallelization friendly

Problem: noticeable recursion overhead

- SNARK generation at each recursion step
- **Concretely** expensive SNARK verifier circuit

Folding Schemes [KST21,BCLMS20,KS23,BC23]

Committed NP Relation:

com: A commitment scheme

$$(x, w) \in R \quad \longrightarrow \quad \begin{array}{l} (x' = (c, x), w) \in R' \\ \text{if and only if} \\ (x, w) \in R \wedge (c = \text{com}(w)) \end{array}$$

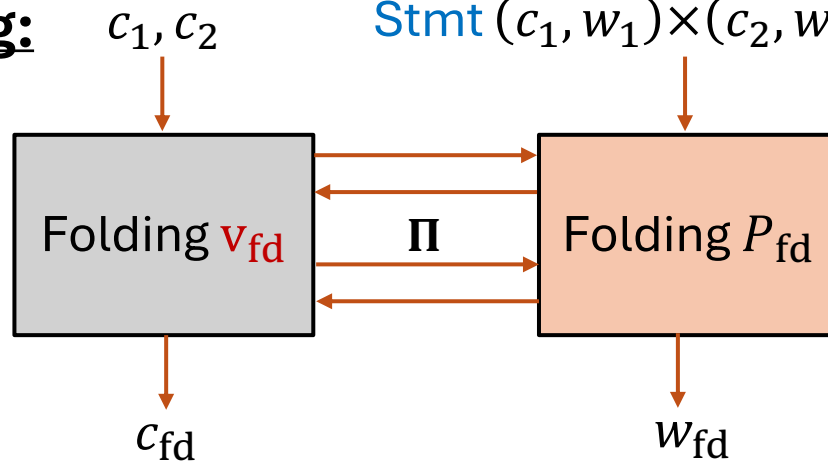
Next: We omit public input x for notational convenience

Folding Schemes [KST21,BCLMS20,KS23,BC23]

E.g., $c_1 = \text{com}(w_1) \wedge \phi(w_1) = 1$

Stmt $(c_1, w_1) \times (c_2, w_2)$

Folding:



Goal: prove $(c_1, w_1) \times (c_2, w_2) \in R \times R$

Reduced goal: prove $(c_{fd}, w_{fd}) \in R$
 P_{fd} and v_{fd} can be made non-interactive

Completeness: If $(c_1, w_1) \times (c_2, w_2) \in R \times R$, then $(c_{fd}, w_{fd}) \in R$ for honest execution

Knowledge soundness: If $(c_{fd}, w_{fd}) \in R$ for P^* 's output w_{fd} , then P^* also knows w_1, w_2

Generalization: Reduction of knowledge [Kothapalli and Parno23]

Input relation: $R_1 := R \times R$

$(c_{in}, w_{in}) := (c_1, w_1) \times (c_2, w_2)$

Output relation: $R_2 := R$

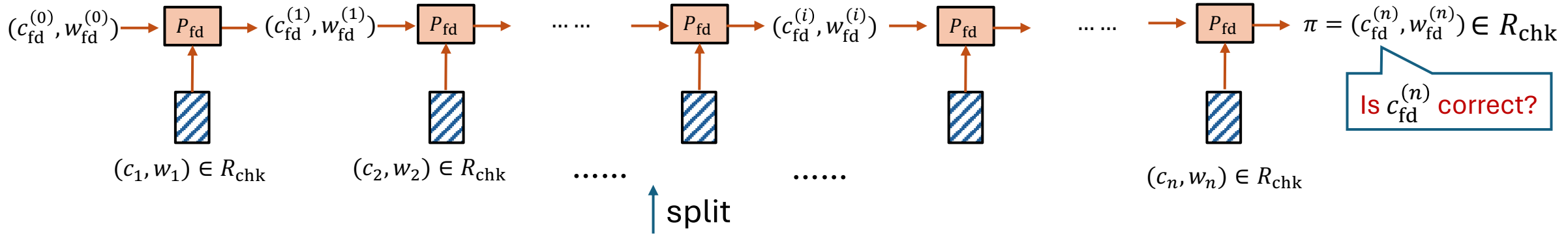
$(c_{out}, w_{out}) := (c_{fd}, w_{fd})$

SNARKs from Folding [KST21,BCLMS20,KS23,B23]

Similar strategies used in SNARGs for P and BARGs [Choudhuru-Jain-Jin21, Waters-Wu22]

Piecemeal SNARK: Prove a chain of computations (can extend to a tree of computations)

SNARK P computes:



Fix:

- Set $x = H(c_n, H(c_{n-1}, \dots, H(c_1)))$ as public input
- SNARK P also sends (c_1, \dots, c_n)
- V checks $x = H(c_n, H(c_{n-1}, \dots, H(c_1)))$ and computes $c_{fd}^{(n)}$ by *iteratively* calling folding v_{fd} given c_1, \dots, c_n

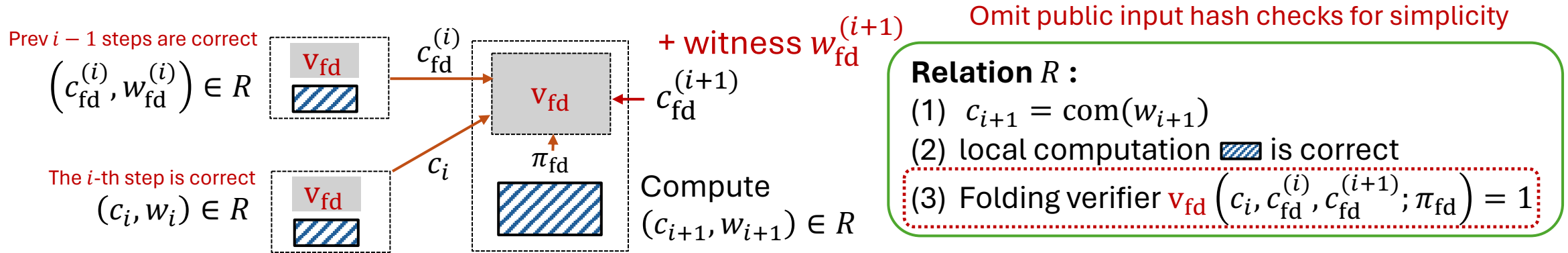
Caveat: proof/verifier complexity linear to n

Idea: Delegate the verifier work into the folded relation

SNARKs from Folding [KST21,BCLMS20,KS23,B^C23]

Piecemeal SNARK: Prove a chain of computations (can extend to a tree of computations)

R : an expanded relation to R_{chk}



Why faster than SNARK recursion? A folding scheme could be more efficient than a SNARK

Folding verifier \mathbf{v}_{fd} :

- \approx check $c_{\text{fd}}^{(i+1)} = c_i + r \cdot c_{\text{fd}}^{(i)}$ for some scalar r
- much simpler than a **SNARK verifier**!

Simpler relation R

Folding prover \mathbf{P}_{fd} :

- $w_{\text{fd}}^{(i+1)} = w_i + r \cdot w_{\text{fd}}^{(i)}$: linear combination of field elems
- much faster than a **SNARK prover**!

Faster folding for relation R than SNARK proving

Folding Schemes: State-of-the-Art

Committed NP statement $(c, w) \in R$

- Instance c : a short com(w) to witness w
- com is linearly-homomorphic for easy folding
e.g., $\text{com}(a) + \text{com}(b) = \text{com}(a + b)$

State-of-the-art:

- Pedersen commitments
 - Linearly-homomorphic
 - Pairing-free
 - No trusted setup

Alternative Option:

Recursive SNARKs from hash-based STARKs

Less efficient: need full SNARK recursion

Security:

- Based on DLOG assumptions & not *post-quantum secure*

Efficiency:

- Require cycle curves
- Prover: many group-exponentiations over a large field
 - Wasteful as real data units usually small (e.g. 32-bit)
- The folding verifier circuit V_{fd} :
 - Elliptic curve scalar multiplications : (
 - Non-native field-op simulations : (
implement arithmetic in \mathbb{F}_q as a circuit over \mathbb{F}_p

Can we construct a folding scheme with

- Post-quantum security
- Ultra-fast prover
- Efficient verifier circuit (e.g., no need for non-native field emulation)

Contributions

LatticeFold: The *first* **lattice-based** folding scheme


- Based on the Module Short-Integer-Solution (**MSIS**) assumption
- Competitive efficiency vs existing folding schemes
 - Linear-time prover + succinct verifier circuit
 - Relatively small fields (e.g., 32-bit or 64-bit)
- Native simulation of ring operations in circuits
 - More friendly for applications like verifiable FHEs/MLs

Technical contribution:

New folding techniques for lattice-based commitments

Warmup: Folding for

Relation R :

- (1) $c_{i+1} = \text{com}(w_{i+1})$
- (2) local computation  is correct
- (3) Folding verifier $\mathbf{v}_{\text{fd}}(c_1, c_2, c_{\text{fd}}; \pi_{\text{fd}}) = 1$

Commitments Opening Relation

Folding for Ajtai Commitment Openings

Committed NP statement $(c, w) \in R$

- Instance c : a short com(w) to witness w
- com is linearly-homomorphic for easy folding

speed \approx Poseidon hash over fast fields [GKRRS19]

How about Ajtai **binding** commitments? [Ajt96,99]

$$\lambda \left\{ \underbrace{A \leftarrow_{\$} \mathbb{Z}_q^{\lambda \times n}}_n \right\} \begin{bmatrix} \vec{w} \end{bmatrix} = \begin{bmatrix} c \end{bmatrix} \in \mathbb{Z}_q^{\lambda} \quad \textbf{Compact}$$

$\vec{w} \in \mathbb{Z}_q^n$
 $\vec{w}_i \in (-\beta, \beta)$ for $i \in [n]$

Binding for “small-norm” \vec{w} (under ~~SIS~~ assumption)

Generalization

$$\mathbb{Z} \Rightarrow R := \mathbb{Z}[X]/(X^d + 1)$$

$$\mathbb{Z}_q \Rightarrow R_q := R/qR$$

[LM'07,PR'07]

Module-SIS

Homomorphic property: (over small-norm messages)

$$\begin{bmatrix} c_1 \end{bmatrix} + \begin{bmatrix} c_2 \end{bmatrix} = \begin{bmatrix} A \end{bmatrix} \times \left(\begin{bmatrix} \vec{w}_1 \end{bmatrix} + \begin{bmatrix} \vec{w}_2 \end{bmatrix} \right) = \begin{bmatrix} A \end{bmatrix} \begin{bmatrix} \vec{w}_1 + \vec{w}_2 \end{bmatrix}$$

How to commit to \vec{w} w/ large norms?

Dealing with Arbitrary Witness

How to commit to an arbitrary witness \vec{w} w/ large norms?

Comm open relation: Our full-fledged protocol fold a similar relation

$$\tilde{R}_{\text{ajtai}}^\beta := \{(c; (\vec{w}, \vec{v})) : (c = A\vec{v}) \wedge (\|\vec{v}\| < \beta) \wedge (\vec{w} = G \times \vec{v})\}$$

Gadget matrix

E.g. $w_1 = [1, 2, 2^2, \dots, 2^{k-1}] \times$

$$\begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \\ \cdot \\ \cdot \\ \vec{v}_k \end{bmatrix}$$

Next, assume that \vec{w} is always **low-norm** in the first place!

Comm open relation:

$$R_{\text{ajtai}}^\beta := \{(c, \vec{w}) : c = A\vec{w} \wedge \|\vec{w}\| < \beta\}$$

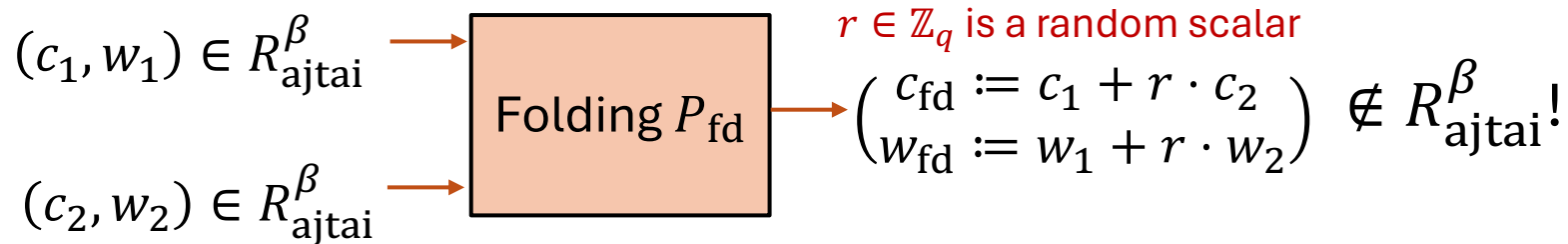
The infinite norm of $w \in \mathbb{Z}^n$
 $\|w\| := \max(|w_i|)_{i=1}^n$

Folding for Ajtai Commitment Openings

Comm open relation: $R_{\text{ajtai}}^\beta := \{(c, \vec{w}) : c = A\vec{w} \wedge \|\vec{w}\| < \beta\}$

The infinite norm of $w \in \mathbb{Z}^n$
 $\|w\| := \max(|w_i|)_{i=1}^n$

Naïve approach:



Problems:

- $\|w_{\text{fd}}\|$ can be larger than β (even if $\|r\|$ is small)
- c_{fd} no longer binding after $\|w_{\text{fd}}\|$ exceeds threshold

Can't support many folding steps

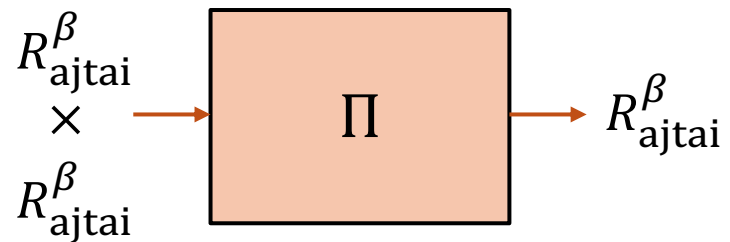
Thoughts:

Make $\|w_1\|, \|w_2\|$ smaller
before random LinComb?

Our Strategy

Relation: $R_{\text{ajtai}}^\beta := \{(c, \vec{w}) : c = A\vec{w} \wedge \|\vec{w}\| < \beta\}$

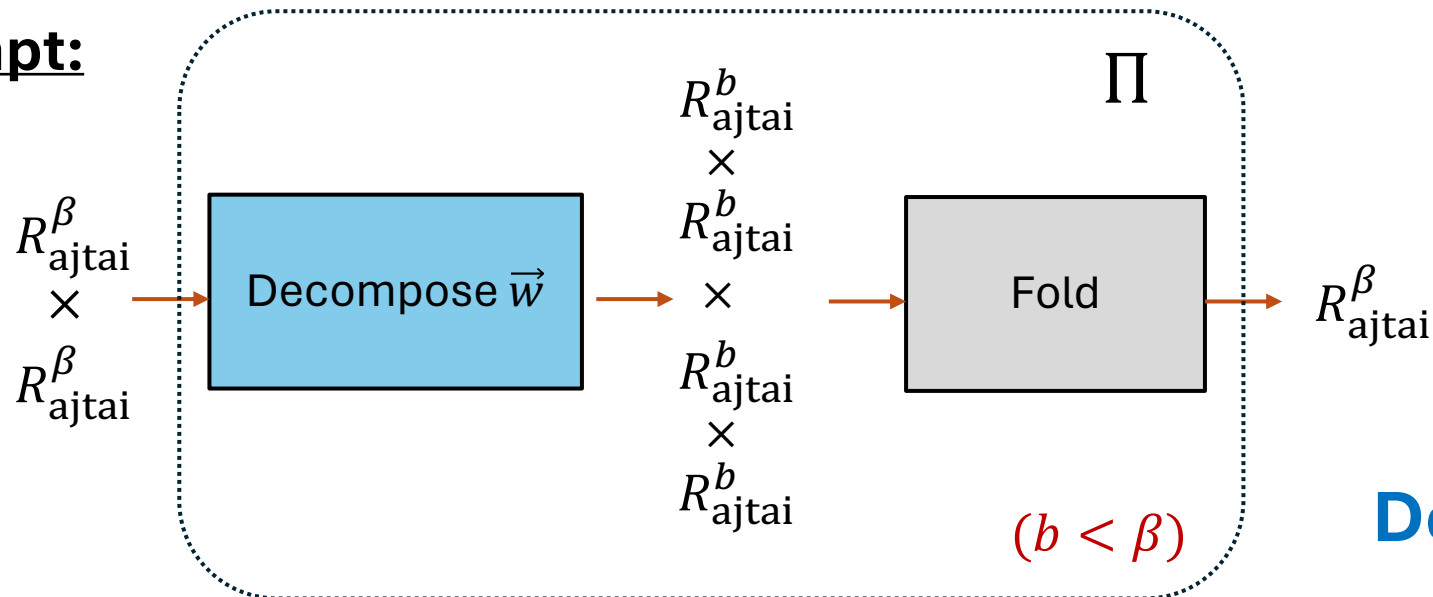
Recall our goal: reduction of knowledge Π Nice property of RoK! [Kothapalli and Parno23]



Sequential composition:

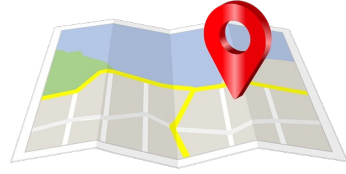


Attempt:

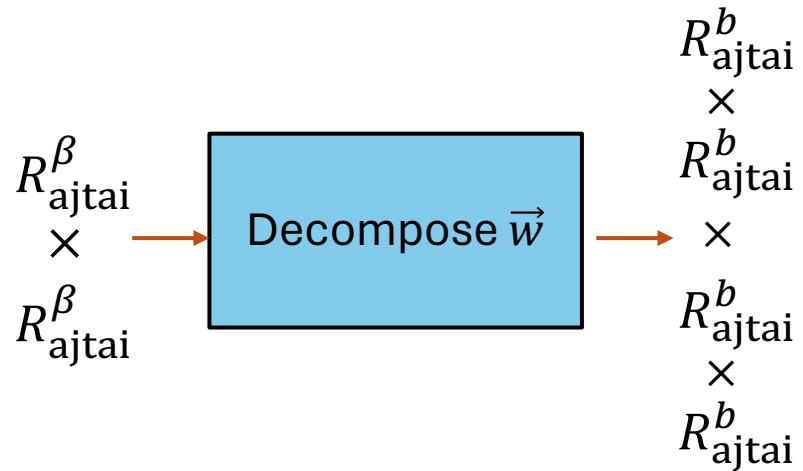


How to instantiate
Decompose and Fold?

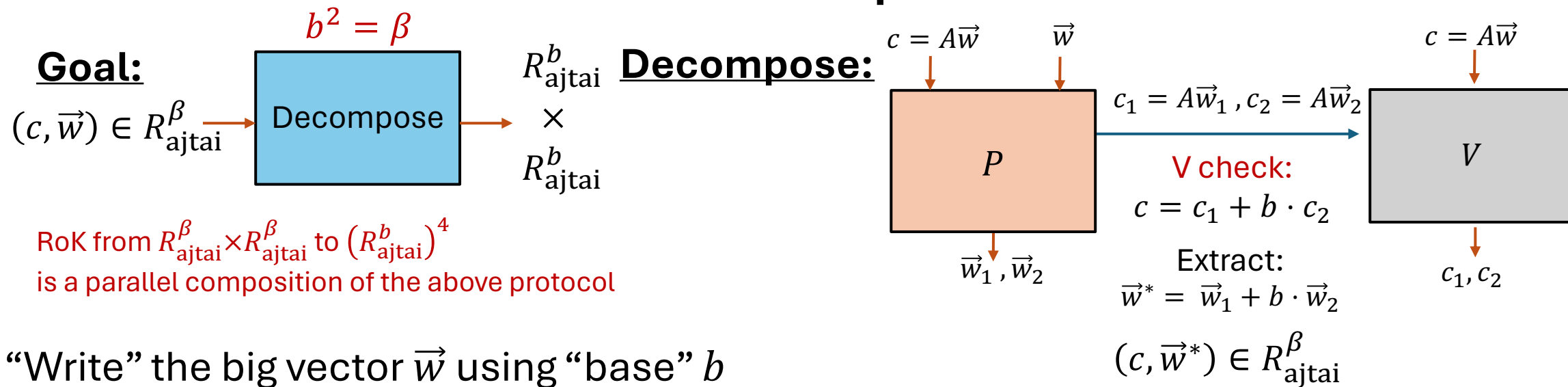
Roadmap



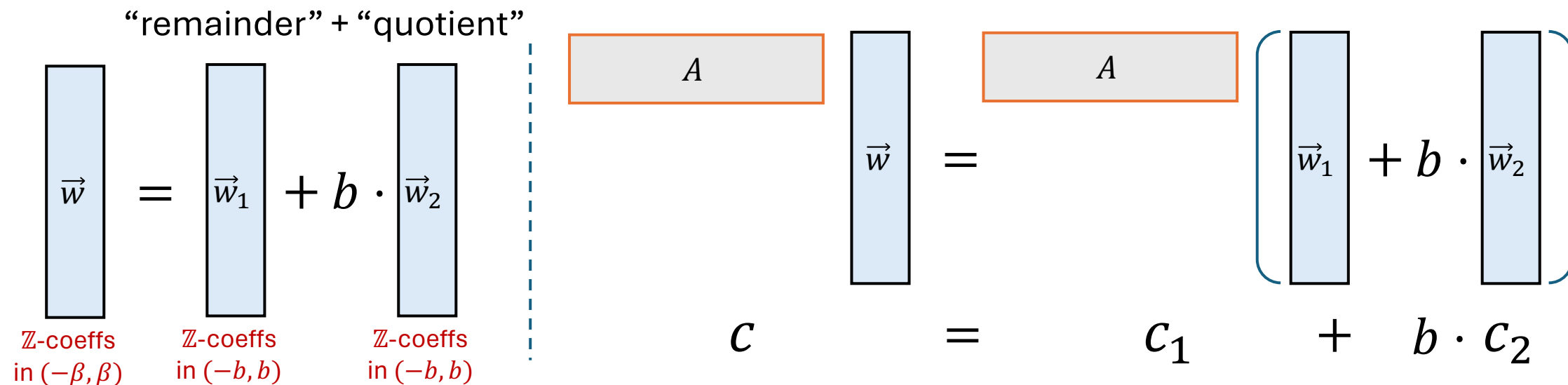
- Decomposition Protocol
- Fold Protocol



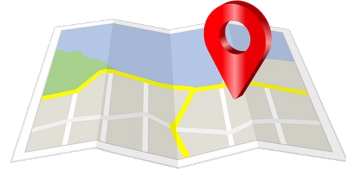
Norm Control with Decomposition



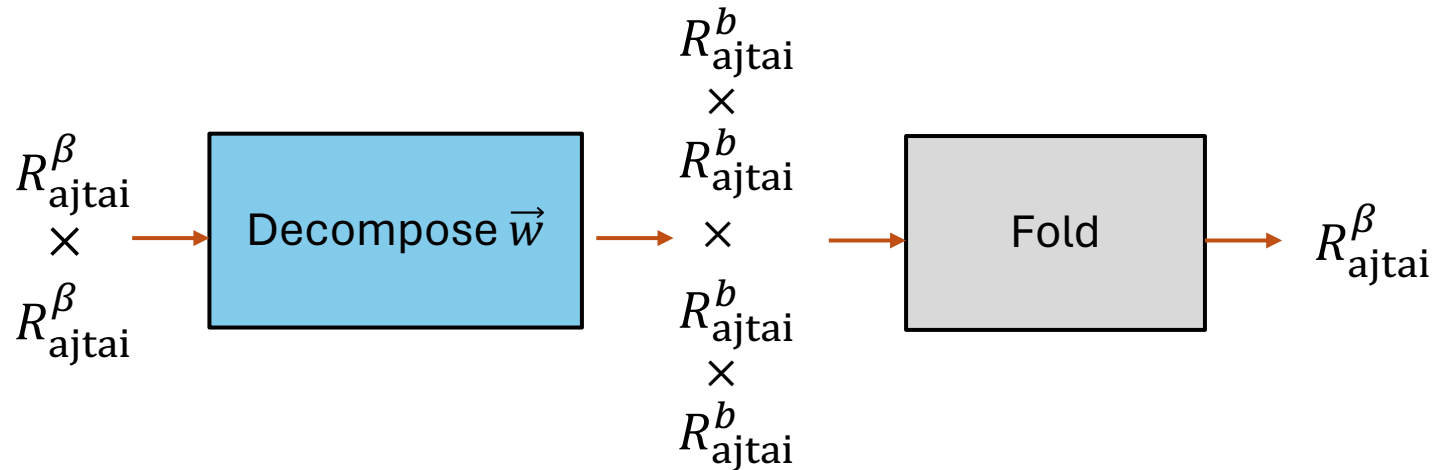
“Write” the big vector \vec{w} using “base” b



Roadmap

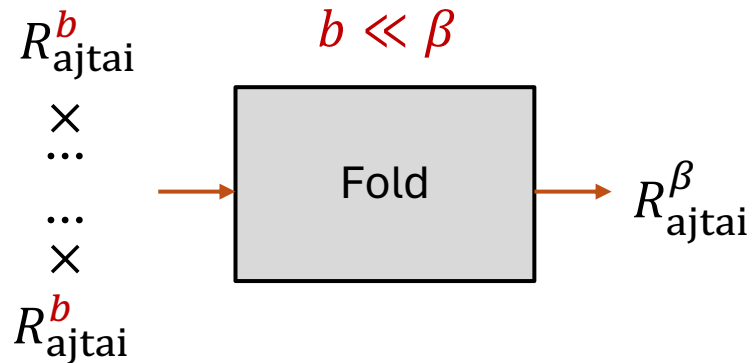


- Decomposition Protocol
- Fold Protocol

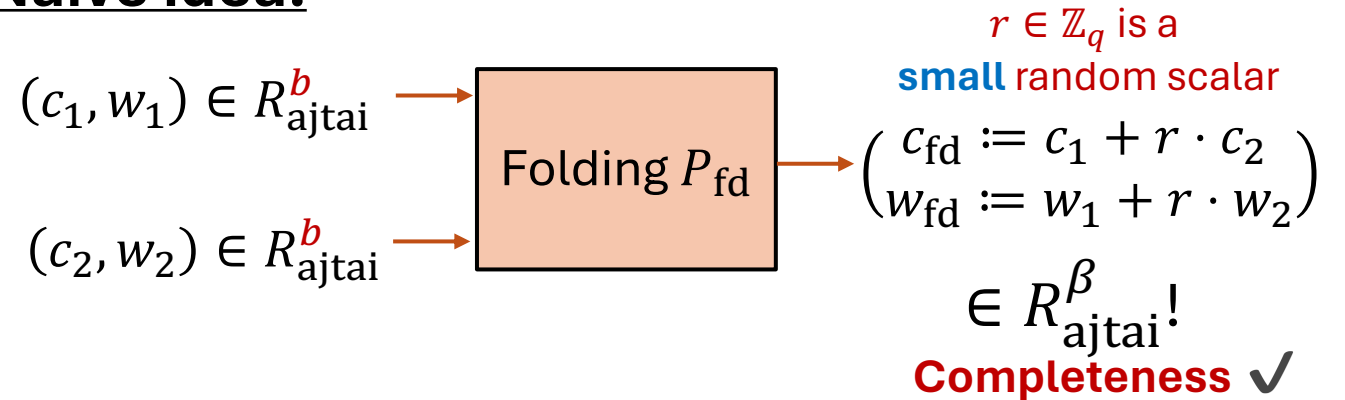


Folding: Naïve Approach

Goal:



Naïve idea:



Knowledge extraction: Rewind P_{fd}^* to obtain $w_{\text{fd}}^x, w_{\text{fd}}^y$ for $c_{\text{fd}}^x = c_1 + r_x \cdot c_2$ and $c_{\text{fd}}^y = c_1 + r_y \cdot c_2$

System of equations:

$$\begin{pmatrix} w_{\text{fd}}^x = w_1 + r_x \cdot w_2 \\ w_{\text{fd}}^y = w_1 + r_y \cdot w_2 \end{pmatrix}$$

Solve linear eqs for w_1, w_2

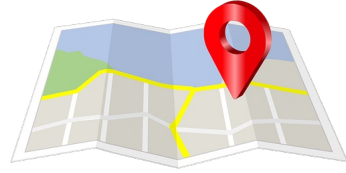
Extracted witness: $(c_2, w_2) \notin R_{\text{ajtai}}^b!$

$$w_2 = \underbrace{(w_{\text{fd}}^y - w_{\text{fd}}^x)}_{\text{The norm can be much larger than } b!} \cdot \underbrace{(r_y - r_x)^{-1}}_{\text{Same for } w_1}$$

The **norm** can be much larger than **b** !

Same for w_1

Roadmap



- Decomposition Protocol

- Fold Protocol

- Naïve extraction + argue **smallness** of the extracted witness

Using Range proof: witness $w \in [-b, b]^n$

(Batched) Range proof via Sumcheck

Goal: Given **input** commitment c' ^{$c' = c_1 \& c_2$ in folding}, prove knowledge of $\vec{w}' = (f_1, f_2, \dots, f_n) \in \mathbb{Z}^n$ ^{$\vec{w}' = w_1 \& w_2$ in folding}

- $c' = A\vec{w}'$
- $\vec{w}' = (f_1, f_2, \dots, f_n)$ **has norm smaller than** b
- Efficient (folding) verifier circuit

Our strategy: Combine naïve folding & extraction + Range proof protocol

- $c' = A\vec{w}'$ (achieved by naïve folding + extraction)
- $\vec{w}' = (f_1, f_2, \dots, f_n)$ **has small norms**

Our solution: A range-proof protocol from Sumcheck

Review of the Sumcheck Protocol [LFKN92]

Goal: Given a “committed” m -variate poly $g(x_1, \dots, x_m)$, convince V that $\sum_{x \in \{0,1\}^m} g(\vec{x}) = s$

Naïve verifier: query g at every $x \in \{0,1\}^m$ and check the sum $\Omega(2^m)$ complexity :

Sumcheck protocol [LFKN92]

- m -round interactive protocol between P and V
 - V sends a random challenge $r_i \in \mathbb{F}$ in each round
- At the end of the protocol, V queries g at a *single* random point

Sumcheck: $\sum_{\vec{x} \in \{0,1\}^m} g(\vec{x}) = s$



Sumcheck protocol [LFKN92] $O(m)$ -time verifier

A reduction from
Sumcheck to Eval stmt

EvalCheck: $g(\vec{r}_1, \dots, \vec{r}_m) = t'$ at a random $\vec{r} \in \mathbb{Z}_q^m$

History: Key ingredient for proving $PH \subseteq IP$ and inspires the proof of $IP = PSPACE$

Goal: Given **input** commitment c' , prove knowledge of $\vec{w}' = (f_1, f_2, \dots, f_n) \in \mathbb{Z}^n$

- $\vec{w}' = (f_1, f_2, \dots, f_n)$ **has norm smaller than b**

Our solution: A range-proof protocol from Sumcheck

Step 1: Rephrase the range-proof statement as a Sumcheck statement

Step 2: Construct a folding protocol for the Sumcheck statement

Step 1: Reducing Range proof to Sumcheck

Range proof: Prove knowledge of a witness $\vec{w}' = (f_1, f_2, \dots, f_n) \in \mathbb{Z}^n$ s.t. Can extend to elements in ring $R = \mathbb{Z}[X]/(X^d + 1)$

$f_1 \in \mathbb{Z}$	f_2	f_{n-1}	f_n
$\in (-b, b) \subseteq \mathbb{Z}$	$\in (-b, b)$				$\in (-b, b)$	$\in (-b, b)$

$$h(x) = 0 \Leftrightarrow x \in (-b, b)$$

$h(x) := x(x+1) \cdot (x-1) \cdots (x+(b-1))(x-(b-1))$ over $\mathbb{Z}_q := \left[-\frac{q}{2}, \frac{q}{2}\right]$ and $q > 2b$ is a prime

$h(f_1) = 0$	$h(f_2) = 0$	$h(f_{n-1}) = 0$	$h(f_n) = 0$
--------------	--------------	-----	-----	-----	------------------	--------------



Embed \vec{w}' to the Boolean hypercube of a multilinear polynomial $f(x_1, \dots, x_{\log n})$

\vec{x}	00 ... 00	00 ... 01	11 ... 10	11 ... 11
$h(f(\vec{x}))$	$h(f_1) = 0$	$h(f_2) = 0$	$h(f_{n-1}) = 0$	$h(f_n) = 0$



Zero-check to sum-check [CBBZ23, Setty20]

Sumcheck: prove that $\sum_{\vec{x} \in \{0,1\}^{\log n}} g(\vec{x}) = 0$ where $g(\vec{x}) := h(f(\vec{x})) \cdot eq_\alpha(\vec{x})$ for a rand $\alpha \in \mathbb{Z}_q^{\log n}$

Step 2: Sumcheck Folding

Range proof: witness $\vec{w}' = (f_1, f_2, \dots, f_n) \in [-b, b]^n$



$$g(\vec{x}) := h(f(\vec{x})) \cdot eq_\alpha(\vec{x})$$

Sumcheck: $\sum_{\vec{x} \in \{0,1\}^{\log n}} g(\vec{x}) = 0$



Sumcheck protocol [LFKN92]

Prover time: $\approx O(bn)$

Verifier time: $O(b \log n)$

EvalCheck: $g(\vec{r}) = t'$ at a random $\vec{r} \in \mathbb{Z}_q^{\log n}$

EvalCheck: $f(\vec{r}) = t$ (and verifier can check $g(\vec{r}) = h(t) \cdot eq_\alpha(\vec{r}) = t'$ itself)

Step 1 ✓

Problem: How to check $f(\vec{r}) = t$ given the comm of f ?

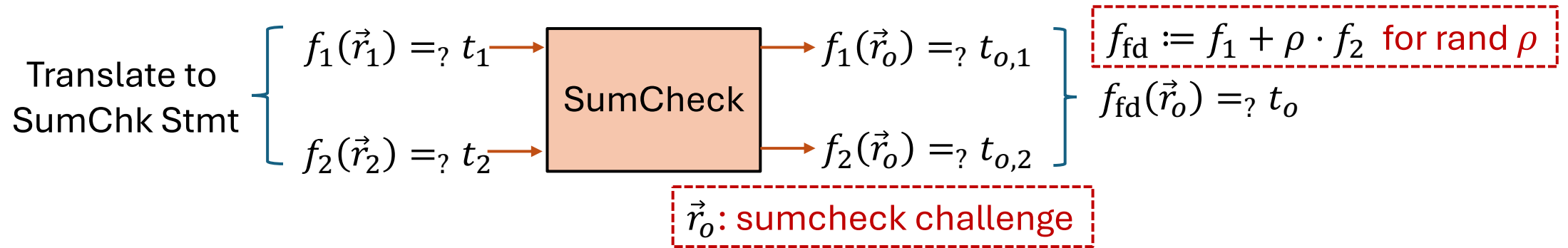
- Send (f_1, f_2, \dots, f_n) to the folding verifier to check it? $O(n)$ folding verifier : (

Observation: EvalStmt $f(\vec{r}) = t$ is easy to fold!

Folding Evaluation Statements

Observation: $f(\vec{r}) = t$ is easy to fold!

Multilinear extension: $f(\vec{r}) = \sum_{\vec{x} \in \{0,1\}^{\log n}} f(\vec{x}) \cdot eq_{\vec{r}}(\vec{x})$ efficiently computable

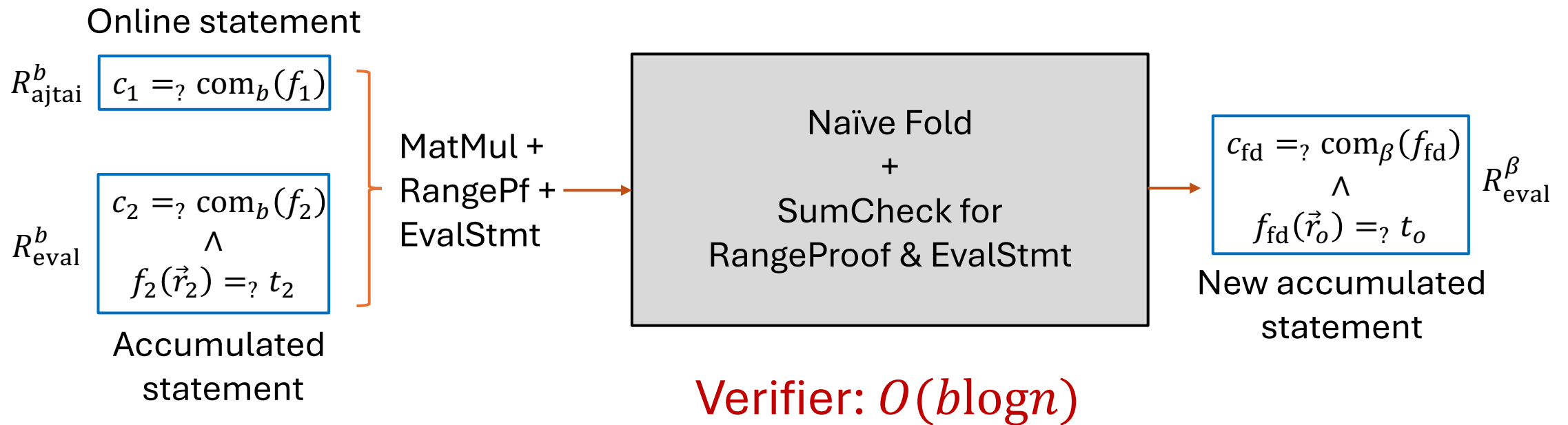


How does it help to check $f(\vec{r}) = t$ given the comm of f ?

- Fold the evaluation statement without checking!

Folding for Ajtai Commitment Openings

Solution: Expand relation R_{ajtai} to include the evaluation statement
 $(c = \text{com}(f)) \wedge (f(\vec{r}) = t)$



The knowledge soundness proof is more subtle than intuition

- A malicious prover can *adaptively* choose the output witness after seeing the challenges
- \Rightarrow The extracted input witnesses could *depend on* the sumcheck challenges

Subtleties & Optimizations

Sumcheck over Rings: [CCKP19, BCS21]

- Ajtai commitments over ring $R_q := \mathbb{Z}_q[X]/(X^d + 1)$ for concrete efficiency
- Small-norm random folding scalar chosen from $S \subseteq R_q$ for negligible soundness error
- Implication: Run Sumcheck over rings

Supporting Small Modulus:

- We want a small modulus q for better efficiency
 - Efficient CPU/GPU ops; no big-number arithmetics
 - More efficient packing of real-world data

Folding for NP-complete relation

Arithmetic over a ring → Great fit for Verifiable ML/FHE

Relation R :

(1) $c_{i+1} = \text{com}(w_{i+1})$

(2) local computation  is correct

(3) Folding verifier $\mathbf{v}_{\text{fd}}(c_1, c_2, c_{\text{fd}}; \pi_{\text{fd}}) = 1$

needs to express computation

Efficiency Estimates

$R_q := \mathbb{Z}_q[X]/(X^{64} + 1) \cong \mathbb{F}_{q^4}^{16}$; q : a 64-bit prime
 C_{chk} : chunk circuit size (e.g. 2^{20} gates over \mathbb{F}_{q^4})
Norm bound: $\beta \approx 2^{16}$; Base: $b = 2$

LatticeFold

speed \approx fast hash

Folding prover: Compute Ajtai commitments
 $O(|C_{\text{chk}}|)$ multiplications over R_q
Can reuse fast FHE impl!

Folding verifier: Sumcheck verifier
 $O(b \cdot \log|C_{\text{chk}}|)$ hashes and R_q -ops
native-ops in the circuit over R_q
Competitive circuit sizes

Existing schemes

Pedersen commitments
 $O(|C_{\text{chk}}|)$ -sized Multi-Scalar-Muls

ECC scalar-mul + (Sumcheck V)
non-native field ops in the circuit
i.e., arithmetic in \mathbb{F}_q as a circuit over \mathbb{F}_p

Piecemeal SNARK proof: ≈ 2 folding instance-witness pairs What if it's still large?

E.g., splitting a stmt of size 2^{40} to 2^{20} chunks \rightarrow 2^{20} -sized chunk stmts

Solution: Use a PQ-secure STARK to prove the correctness of the folding statement

$< 100\text{KB}$ and 2ms verifier (STIR[ACFY24])

$< 5\text{KB}$ w/ Hyperplonk+KZG[CBBZ23]

Summary & Open Problems

Takeaway:

- The *first* lattice-based folding scheme based on Ajtai commitments
- Gives memory-efficient, plausibly PQ-secure SNARKs, with fast provers
- Generic techniques for folding lattice-based commitments w/ norm constraints

Open problems:

- *Compact + homomorphic* lattice commitments with **no norm constraints**
- Folding table lookup relations (e.g., from Lasso [Setty-Thaler-Wahby23])
- Efficient implementation

Concurrent work:[Bünz-Mishra-Nguyen-Wang24]

- Purely from hashing; no lattice crypto
- General optimization techniques for piecemeal SNARKs (apply to LatticeFold)
- Larger verifier circuit; only supports bounded-depth folding (attack exists)

Thank you!

<https://eprint.iacr.org/2024/257.pdf>

Expecting updates soon!