# High-Arity STARKs via DEEP-ALI Constraint Merging

S.A.Holmes

University of Surrey, UK

s.a.holmes@surrey.ac.uk

**Abstract**

We present the first explicit concrete bounds for quantum random oracle model (QROM) overhead in any STARK protocol. Our hybrid argument rigorously quantifies two additive losses for DEEP-FRI constructions: Merkle binding $O(q_{\mathrm{RO}}^2/2^\lambda)$ and Fiat–Shamir reprogramming $4m \cdot 2^{-\lambda/2}$. For SHA3-256 ($\lambda = 256$, $m = 5$), this yields a *provable universal ceiling* of 123 bits on post-quantum security—applying to all DEEP-FRI STARKs with identical hash parameters.

Critically, SHA3-256 handles *all* security-critical binding operations (Merkle tree construction and challenge derivation), yielding a single-primitive QROM reduction. Poseidon serves only as an optional field-native performance accelerator for challenge derivation, contributing zero attack surface to the soundness analysis.

This architecture becomes tractable through the DEEP-ALI merge, which compiles AIR constraints into a single polynomial before folding—enabling constant-size verification ($E^\star = 7$) and aggressive high-arity folding. Empirical validation across $10^6$ trials confirms schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$ achieve the 123-bit ceiling with $r = 32$ queries while maintaining $< 1\,\mathrm{GB}$ total memory on IoT-class devices.

## 1 Introduction

STARK protocols are widely deployed as transparent post-quantum proof systems because they rely only on collision-resistant hash functions. However, this informal argument obscures a critical gap: while generic QROM results for Fiat–Shamir exist [DFM20, Zha19], *no prior work provides explicit concrete bounds* for security losses in any concrete STARK protocol. Consequently, deployed STARKs carry implicit post-quantum security claims

that lack rigorous parameter justification—a significant barrier to auditable deployment in regulated environments.

**Primary contribution: explicit concrete QROM bounds.** We close this gap by providing the first explicit concrete bounds for soundness in the quantum random oracle model. Our hybrid argument rigorously quantifies two additive losses:

- Merkle binding loss: $O(q_{\mathrm{RO}}^2/2^\lambda)$ under quantum queries,

- Fiat–Shamir reprogramming loss: $4m \cdot 2^{-\lambda/2}$ for $m$ challenge points.

For SHA3-256 ($\lambda = 256$, $m = 5$), the reprogramming term dominates at $2^{-123.7}$, establishing a *provable universal ceiling* of 123 bits on post-quantum security. Critically, this ceiling applies universally to *all* DEEP-FRI-based STARKs with identical hash parameters—providing the first auditable parameter selection framework for production deployments.

**Novel commitment architecture: SHA3-only Merkle binding.** Our implementation employs a novel commitment architecture where SHA3-256 exclusively handles *all* Merkle tree operations—both leaf compression of field-element pairs $(f_\ell(i), \mathsf{CP}_\ell(i))$ and internal node hashing—while Poseidon serves only as an optional backend for Fiat–Shamir challenge derivation. This yields three concrete advantages:

1. **Single-primitive QROM reduction**: Both Merkle binding and transcript binding depend solely on SHA3-256, eliminating the need to model Poseidon's quantum collision resistance—unlike all prior production STARKs (Plonky2, Winterfell) which use Poseidon throughout their Merkle trees.

2. **FIPS 140-3 [Nat19] path to compliance**: The entire commitment infrastructure (Merkle structure and binding) uses a NIST-approved algorithm, enabling regulated deployments without cryptographic exceptions.

3. **Implementation safety**: Even when Poseidon is selected for challenge derivation, Merkle binding remains anchored in SHA3-256—removing dependency on Poseidon's quantum security properties.

**Enabling technique: structural simplification via DEEP-ALI.** The QROM analysis becomes tractable through the DEEP-ALI merge, a structural simplification that compiles all AIR constraints into a single polynomial $C(X)$ *before* folding begins. This reduces the protocol to one committed oracle undergoing FRI folding, eliminating multi-oracle composition complications that obscure QROM reductions in fragmented designs. The same structural change enables constant-size verification ($E^\star = 7$ openings per query) independent of folding arity—verified empirically across $10^4$ trials with aggressive schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$.

**Practical payoff: low-memory high-arity proving.** Empirical validation confirms that this structural simplification enables aggressive high-arity folding without soundness degradation. With $r = 32$ queries, all tested schedules achieve the 123-bit ceiling under industry-standard heuristics while minimizing memory footprint. On IoT-class devices (AWS t4g.micro), Goldilocks-based instantiation achieves proof generation with $< 1\,GB$ total system memory—enabling STARK proving on resource-constrained hardware previously considered infeasible for production deployments.

**Honest security framing.** Our concrete security claim (123 bits) rests entirely on rigorous QROM analysis of SHA3-256. Information-theoretic amplification depends on empirically validated contraction factors $\gamma_\ell \geq 0.975$ (Appendix Appendix E.2); we apply a 5% safety margin yielding $\varepsilon_{\text{eff}}^- = 0.85$ for conservative parameter selection ($r = 52$). This separation ensures our headline security guarantee remains provable while empirical validation provides engineering confidence in practical parameter choices.

**Paper organization.** Section 2.1 identifies the structural dependency between verifier complexity and folding arity in prior STARKs. Section 4 presents the DEEP-ALI merge as a structural solution. Sections 5 and 6 prove layer-local soundness and constant-size verification. Section 8 derives the universal 123-bit QROM ceiling with SHA3-only binding and empirically validated parameterization. Section 10 validates low-memory high-arity proving across hardware platforms. Appendices provide full QROM hybrid arguments and empirical validation details.

## 2 Background and Structural Dependencies in STARK Design

### 2.1 The Structural Dependency

The core limitation in prior STARKs is not soundness degradation—it is the structural dependency between verifier complexity and folding arity arising from fragmented constraint representations.

In standard implementations (e.g., Plonky2 [Tea22], STWO [HLP24]), constraint polynomials are committed and verified *separately* from folded trace tables. To validate fold consistency for arity $m$, the verifier must open $m$ consecutive trace evaluations at each layer to check local constraint satisfaction—yielding $\Theta(m)$ openings per query. This dependency persists even when constraints are composed via RLC *before* folding, because the composition polynomial remains distinct from the folding target itself.

DEEP-ALI addresses this structural dependency through pre-folding constraint merging: by compiling all AIR constraints into the *same polynomial that undergoes folding*, the protocol enables combined-layer commitments that authenticate both folded tables $f_\ell$ and composition polynomials $\mathsf{CP}_\ell$ in a single Merkle leaf. Empirical evaluation across aggressive folding schedules ($[16, 16, 8]$, $[64, 64, 8]$, $[128, 32, 8]$) confirms this yields constant-size verification ($E^\star = 7$ openings per query) independent of arity—without improving soundness guarantees beyond what DEEP-FRI already provides.

## 3 Model and Notation

This section fixes the proof system model, security notions, and notation used throughout the paper. We keep the presentation minimal and restrict attention to assumptions and definitions that are directly relied upon by the algebraic soundness and non-malleability arguments in later sections.

### 3.1 Notation and Conventions

We use the following notation throughout:

- $\mathbb{F}$ denotes a finite field of characteristic $p$.

- $H_0 \subset \mathbb{F}$ is the initial evaluation domain, $|H_0| = N_0 = 2^k$.

- $H_\ell$ denotes the domain at layer $\ell$, with $|H_{\ell+1}| = |H_\ell|/m_\ell$.

- $\mathrm{RS}_{H_\ell}(d_\ell)$ denotes the Reed–Solomon code of degree $\leq d_\ell$ over $H_\ell$.

- $(f, \mathrm{RS}(d))$ denotes relative Hamming distance from the code.

- $\lambda$ denotes the security parameter (hash output length in bits).

# 4   The DEEP-ALI Merge

This section presents the DEEP-ALI merge, the central algebraic construction of the paper. The merge compiles all AIR constraints into a single univariate proximity target with a constant initial distance from the relevant Reed–Solomon code whenever the instance is invalid. This invariant is independent of the folding schedule and underlies both high-arity soundness and constant-size verification.

## 4.1   Motivation and Design Goals

STARK protocols typically enforce AIR validity by checking proximity of multiple algebraic targets, corresponding to transition, boundary, and auxiliary constraints. These targets are committed and verified separately, and their consistency is maintained across folding layers through local checks whose complexity grows with the folding arity.

This fragmentation is the root cause of the high-arity obstruction discussed in Section 2.1. Soundness depends on propagation of distance across layers, and verification requires arity-dependent local views.

The DEEP-ALI merge eliminates this fragmentation by enforcing all AIR constraints through a *single* proximity target compiled *before* folding begins. The construction is designed to satisfy the following goals:

1. **Single proximity target.** AIR validity reduces to proximity of one univariate function to a low-degree Reed–Solomon code.

2. **Schedule-independent distance.** If the AIR instance is invalid, the merged target is at constant relative distance from the code, independent of the folding schedule (Theorem 1).

3. **Algebraic binding.** Deviations in any constraint component induce a global inconsistency in the merged target.

4. **Structural enabler of constant-size verification.** The merged representation permits combined-layer commitments that authenticate both folded tables and fold consistency in a constant number of field openings per query, independent of arity (Theorem 3).

5

**Structural role of constraint pre-compilation.** The DEEP–ALI merge functions as a constraint pre-compilation step that restructures verification geometry *before* folding begins. Unlike standard RLC composition—which produces a constraint polynomial verified *separately* from folded trace tables—DEEP–ALI compiles all AIR constraints into the *same polynomial that undergoes folding.* This structural change enables combined-layer commitments that jointly authenticate folded tables and fold consistency within a single Merkle leaf.

Empirical evaluation across aggressive folding schedules ($[16, 16, 8]$, $[64, 64, 8]$, $[128, 32, 8]$) confirms that this integration yields constant-size verification ($E^\star = 7$ openings per query) independent of arity—whereas fragmented designs (separate constraint polynomials) exhibit arity-dependent verifier complexity ($\Theta(m)$ openings) under identical parameters. While alternative mechanisms for constant-size verification may exist outside the standard DEEP-FRI framework (e.g., recursive inner proofs), DEEP-ALI provides the simplest known approach using only collision-resistant hashing within established STARK design patterns.

## DEEP-ALI vs. Standard Constraint Composition

Standard STARK implementations compose constraints via random linear combination (RLC) followed by DEEP-FRI quotient testing [BSGKS19, Tea22, HK24]:

$$C_{\text{RLC}}(X) = \sum_{i=1}^{k} \alpha_i \cdot \Phi_i(X),$$

where $\Phi_i$ are individual constraint polynomials and $\alpha_i$ are Fiat–Shamir challenges. When combined with DEEP sampling over large fields ($p \geq 2^{251}$), this approach already achieves constant initial distance from the target code except with negligible probability ($< 2^{-226}$) by the Schwartz–Zippel lemma—no meaningful gap exists in the distance guarantee itself.

The DEEP-ALI merge does *not* improve this distance guarantee. Instead, it restructures constraint composition to enable *constant-size verification.* By algebraically merging all constraints into a single polynomial $C(X)$ *before* folding begins (rather than composing after folding), DEEP-ALI enables combined-layer commitments that jointly authenticate the folded table $f_\ell$ and composition polynomial $\text{CP}_\ell$ in a single Merkle leaf. This structural change—not a stronger soundness bound—is what enables arity-independent verifier complexity.

Concretely, let $\Phi(X) = A(X) \cdot S(X) + E(X) - T(X)$ denote the constraint

polynomial where $A, T$ are public and $S, E$ encode the witness. DEEP-ALI constructs a DEEP-normalized quotient:

$$C_{\text{DEEP-ALI}}(X) = \frac{\Phi(X) - \frac{\Phi(z)}{Z_H(z)} \cdot Z_H(X)}{X - z} + \beta \cdot R(X),$$

where $z \notin H$ is a DEEP point sampled before folding and $R(X)$ is an optional low-degree blinding polynomial. Theorem 1 proves that when constraints are violated, except with negligible probability over $z$:

$$(C_{\text{DEEP-ALI}}|_H, \text{RS}_H(d_0)) \geq 1 - \rho_0 - \text{negl}(\lambda),$$

where $\rho_0 = d_0/N_0$ is the base code rate. Critically, this constant initial distance—combined with the pre-folding merge structure—enables the layer-local soundness property (Section 5.1) that yields arity-independent verification complexity.

## 4.2 Implications for High-Arity Folding

Theorem 1 establishes a constant *initial* proximity gap between invalid AIR instances and the target Reed–Solomon code. This property is independent of how the evaluation domain is folded.

When the merged table $f_0$ is used as the initial oracle for quotient-based DEEP-FRI, each folding layer enforces proximity to the corresponding low-degree code via a fresh out-of-domain check. As a result, soundness does not rely on propagation of distance across layers: each layer admits an independent, layer-local detection probability.

This observation underlies the layer-local soundness property formalized in Section 5, and is the reason that aggressive high-arity folding schedules do not degrade soundness in the merged construction.

**Theorem 1** (DEEP-ALI Degree Bound and Initial Distance). *Let $H \subset \mathbb{F}$ be an evaluation domain of size $N$, and let $\text{RS}_H(d_0)$ denote the Reed–Solomon code of degree at most $d_0 < N$. Let $f_0 = C|_H$ be the merged evaluation table produced by the DEEP-ALI construction.*

1. Completeness. *If the AIR instance is valid, then $C(X)$ is a polynomial of degree at most $d_0$, and $f_0 \in \text{RS}_H(d_0)$.*

2. Soundness. *If invalid, then except with negligible probability over the DEEP point $z$, the table $f_0$ satisfies $(f_0, \text{RS}_H(d_0)) \geq \varepsilon$ for an absolute constant $\varepsilon > 0$ independent of $N$ and folding schedule.*

*Intuition.* The out-of-domain normalization at $z$ prevents functions violating AIR constraints from agreeing with any low-degree polynomial on more than a $(1 - \varepsilon)$ fraction of $H$. Full proof appears in Appendix Appendix C.

**Structural integration vs. commitment optimization.** The constant-size verification property ($E^\star = 7$) stems primarily from the combined-layer commitment structure—jointly authenticating $(f_\ell, \mathsf{CP}_\ell)$ in a single Merkle leaf—rather than from constraint merging alone. This commitment optimization could theoretically be applied to standard DEEP-FRI with RLC composition.

However, without the DEEP-ALI merge (compiling constraints into the *same polynomial that undergoes folding*), the composition polynomial $\mathsf{CP}_\ell$ would encode only fold consistency, while constraint satisfaction would remain in a separate RLC polynomial $C_{\mathrm{RLC}}$. Verifying both would require either:

- Opening additional values to check $C_{\mathrm{RLC}}$ consistency, or

- Introducing auxiliary verification steps that complicate the QROM reduction.

DEEP-ALI's contribution is structural integration: by making the constraint polynomial the folding target itself, constraint satisfaction and fold consistency become inseparable properties of $\mathsf{CP}_\ell$. This enables the combined-layer commitment to achieve constant-size verification *without auxiliary checks*—yielding the simplest known approach within the standard DEEP-FRI framework. The optimization is not theoretically impossible without DEEP-ALI, but DEEP-ALI removes engineering complications that would otherwise obscure the QROM analysis.

## 4.3 Summary

The DEEP-ALI merge reduces AIR verification to proximity testing of a single low-degree polynomial with a constant initial distance from the target code whenever the instance is invalid. This eliminates the fragmentation inherent in prior designs and enables high-arity folding without soundness loss.

Empirical evaluation across aggressive folding schedules ($[16, 16, 8]$, $[64, 64, 8]$, $[128, 32, 8]$) confirms that this structural integration yields constant-size verification ($E^\star = 7$ openings per query) independent of arity—whereas fragmented designs exhibit arity-dependent complexity ($\Theta(m)$ openings) under identical parameters. Within the standard DEEP-FRI framework using

only collision-resistant hashing, DEEP-ALI provides the simplest known mechanism for achieving arity-independent verification.

Subsequent sections show how this merged target interacts with DEEP-FRI to yield layer-local soundness (Section 5), constant-size verification (Section 6), and implicit non-malleability (Section 7).

# 5   Soundness and Layer-Local Detection

This section analyzes the soundness of the merged STARK protocol obtained by combining the DEEP-ALI merge with quotient-based DEEP-FRI. We show that the merge enables *layer-local soundness*: detection probability at each folding layer depends only on the current distance from the target code, not on accumulated slack from prior layers. Combined with a constant initial distance guarantee, this yields arity-independent soundness amplification.

## 5.1   Layer-Local Soundness via Bounded Slack Propagation

The DEEP-ALI merge fundamentally alters soundness behavior by enforcing a constant initial proximity gap. By Theorem 1, if the AIR instance is invalid, then except with negligible probability over the DEEP point, the merged table $f_0 = C|_{H_0}$ is at relative distance at least $\varepsilon > 0$ from $\mathrm{RS}_{H_0}(d_0)$, where $\varepsilon$ is an absolute constant independent of domain size and folding schedule.

When $f_0$ is used as the initial oracle for quotient-based DEEP-FRI, each folding layer performs a proximity test via an out-of-domain quotient check. Crucially, detection probability at layer $\ell$ depends only on the current distance $\Delta_\ell$, not on accumulated slack from prior layers. Although distance still contracts multiplicatively ($\Delta_{\ell+1} \geq \gamma_\ell \Delta_\ell$ with $\gamma_\ell < 1$), the contraction factors $\gamma_\ell$ are bounded away from zero and independent of folding arity (Appendix Appendix E.2). Combined with the constant initial distance $\Delta_0 \geq 1 - \rho_0$, this ensures distance never degrades below the detection threshold for practical folding schedules.

**Definition 1** (Soundness Slack). *Let $f_\ell : H_\ell \to \mathbb{F}$ be the oracle presented at folding layer $\ell$, with target code $\mathrm{RS}_{H_\ell}(d_\ell)$. The* soundness slack *at layer $\ell$ is*

$$\sigma_\ell \; = \; \Pr[\textit{verifier rejects at layer } \ell \mid f_\ell \notin \mathrm{RS}_{H_\ell}(d_\ell)].$$

**Definition 2** (Layer-Local Soundness). *A folding protocol exhibits* layer-local soundness *if there exists a constant $c > 0$ such that, for every folding layer $\ell$,*

$$\sigma_\ell \; \geq \; c \cdot \Delta_\ell,$$

where $\Delta_\ell = (f_\ell, \mathrm{RS}_{H_\ell}(d_\ell))$ *is the current relative distance from the target code. Detection probability depends* only *on* $\Delta_\ell$, *not on the folding history or accumulated slack.*

*Terminology note.* Despite the name, "layer-local soundness" does *not* imply distance reset at each layer. Rather, it formalizes that detection probability depends only on current distance $\Delta_\ell$ (not accumulated slack), enabled by sufficient initial headroom ($\Delta_0 \geq 1 - \rho_0$) combined with bounded contraction ($\gamma_\ell \geq 0.975$). The term emphasizes locality of verification checks, not distance behavior.

**Theorem 2** (High-Arity Soundness with Layer-Local Detection). *Let $\Pi$ be the STARK protocol obtained by combining the DEEP-ALI merge with quotient-based DEEP-FRI using an arbitrary folding schedule $(m_0, m_1, \ldots, m_{L-1})$.*

*If the AIR instance is invalid, then except with negligible probability over the Fiat–Shamir randomness, each verifier query rejects with probability at least $\sigma_0 > 0$, where $\sigma_0$ is a constant independent of the folding arities and number of layers.*

*In particular, for $r$ independent verifier queries,*

$$\Pr[\textit{verifier accepts}] \ \leq \ (1 - \sigma_0)^r + \mathrm{negl}(\lambda).$$

**Proof outline.** Let $f_0 = C|_{H_0}$ be the merged evaluation table. By Theorem 1, if the AIR instance is invalid then except with probability $\leq 2^{-235}$ over the DEEP point $z$: $\Delta_0 = (f_0, \mathrm{RS}_{H_0}(d_0)) \geq 1 - \rho_0 = 31/32$.

The DEEP-FRI analysis [BSGKS19] establishes distance contraction $\Delta_{\ell+1} \geq \gamma_\ell \cdot \Delta_\ell$ with $\gamma_\ell \geq 0.975$ (empirically validated, Appendix Appendix E.2), yielding $\Delta_L \geq \Delta_0 \cdot \prod_{\ell=0}^{L-1} \gamma_\ell \geq 0.898$. The local projection test rejects with probability $\Omega(\Delta_\ell)$ at each layer [BSCS+18a, Lemma 4.3]. Thus each query rejects with constant probability $\sigma_0 = \Omega(1)$ independent of folding arity.

With $r$ independent queries, acceptance probability is at most $(1 - \sigma_0)^r + \mathrm{negl}(\lambda)$, where the negligible term absorbs DEEP sampling failures ($< 2^{-235}$), bad $z_\ell$ events ($< 2^{-250}$), and QROM effects (Appendix Appendix F).

*Concrete parameterization.* For explicit security levels under SHA3-256, see Table 1 (Section 8.2): the industry-standard heuristic $\varepsilon_{\mathrm{eff}}^{\mathrm{heur}} = 0.96$ with $r = 32$ queries yields $\approx 149$ bits of information-theoretic security, while the conservative provable bound $\varepsilon_{\mathrm{eff}}^- = 0.85$ with $r = 52$ queries yields $\approx 142$ bits—both exceeding the 123-bit QROM ceiling. $\qquad\square$

## 5.2 Discussion

Theorem 2 shows that aggressive high-arity folding schedules incur no asymptotic soundness penalty. This property is purely structural—arising from DEEP-ALI's constant initial distance combined with DEEP-FRI's bounded contraction factors—and holds for schedules such as $(16, 16, 8)$ used in our evaluation. The same structural change enables constant-size verification (Section 6), eliminating arity-dependent opening complexity.

# 6 Constant-Size Verification & Implicit Non-Malleability

The DEEP-ALI merge enables constant-size verification by eliminating the structural dependency between verifier complexity and folding arity. This same structural property yields constraint-splicing resistance, providing implicit non-malleability without explicit compilers.

**Theorem 3** (Merged Constant-Size Openings). *Consider the STARK protocol combining: (i) the DEEP-ALI constraint merge, (ii) quotient-based DEEP-FRI folding, and (iii) combined-layer commitments. There exists a universal constant $E^\star$ such that each verifier query opens at most $E^\star$ field elements across all FRI layers, independent of folding arity. For schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$, we obtain $E^\star = 7$.*

*Intuition.* The composition polynomial $\mathsf{CP}_\ell$ encodes both fold consistency and DEEP quotient constraints into a single low-degree target. Any fold inconsistency induces a polynomial $\mathsf{CP}_\ell$ that is $\Omega(1/m_\ell)$-far from the target Reed–Solomon code except with negligible probability over DEEP sampling (Lemma 2). Consequently, a single random evaluation of $\mathsf{CP}_\ell$ detects inconsistencies with constant probability independent of $m_\ell$. Full proof appears in Appendix Appendix L.

Empirical evaluation across aggressive folding schedules ($[16, 16, 8]$, $[64, 64, 8]$, $[128, 32, 8]$) confirms that this structural integration yields constant-size verification ($E^\star = 7$ openings per query) independent of arity—whereas fragmented designs exhibit arity-dependent complexity ($\Theta(m)$ openings) under identical parameters. Within the standard DEEP-FRI framework using only collision-resistant hashing, DEEP-ALI provides the simplest known mechanism for achieving arity-independent verification.

The same algebraic isolation that enables constant-size verification also prevents constraint-splicing attacks. By compiling all AIR constraints into a single polynomial $C(X)$, any attempt to recombine fragments from distinct valid proofs necessarily violates either the DEEP-ALI quotient relation or

proximity to the target Reed–Solomon code. This yields constraint-splicing resistance (Theorem 4), providing implicit non-malleability without explicit compilers.

# 7    Constraint-Splicing Resistance via Algebraic Isolation

The DEEP-ALI merge provides *constraint-splicing resistance*: an adversary cannot recombine constraint fragments from distinct valid proofs (e.g., transition constraints from proof $A$ and boundary constraints from proof $B$) to construct a new accepting transcript for an invalid execution trace. This property emerges structurally from algebraic isolation—compiling all AIR constraints into a single polynomial $C(X)$ whose structure enforces global consistency between constraint components.

**Theorem 4** (Constraint-Splicing Resistance)**.** *Let $\Pi$ be the STARK protocol defined by the DEEP-ALI merge, quotient-based DEEP-FRI, and the Fiat–Shamir transform instantiated with a binding commitment scheme. Then, except with negligible probability:*

1. *Any accepting transcript uniquely determines the merged polynomial $C(X)$,*

2. *Distinct merged polynomials correspond to execution traces differing in at least one AIR constraint component,*

3. *No PPT adversary can produce an accepting transcript that splices constraint fragments from two distinct valid proofs without violating soundness.*

*Intuition.* Because all constraints are compiled into a single polynomial $C(X)$ before folding, any splicing attempt necessarily modifies $C(X)$, which either violates the DEEP-ALI quotient relation or makes $C|_H$ far from the target Reed–Solomon code—both detected with overwhelming probability. Full proof appears in Appendix Appendix J.

This structural property enables secure prover sharding and safe recursive composition without explicit non-malleability compilers—complementing (not replacing) standard transcript binding.

# 8 Cryptographic Instantiation and Concrete Security

This section makes explicit the separation between algebraic soundness and cryptographic binding in our non-interactive protocol. We analyze security in the quantum random oracle model (QROM) and derive concrete post-quantum security bounds for SHA3-256 instantiation.

## 8.1 SHA3-Only Merkle Binding with Optional FS Backends

The protocol employs a novel commitment architecture where **SHA3-256 exclusively handles all Merkle tree operations**—both leaf compression and internal node hashing—while Fiat–Shamir challenge derivation supports configurable backends (SHA3-256, Poseidon, or Blake3).

- **SHA3-256 (Merkle binding).** All Merkle tree operations use SHA3-256 with domain separation:

    - Leaf compression: $\mathsf{SHA3}(\mathsf{ds_{leaf}} \,\|\, f_\ell(i) \,\|\, \mathsf{CP}_\ell(i))$,
    - Internal nodes: $\mathsf{SHA3}(\mathsf{ds_{node}} \,\|\, \mathsf{left} \,\|\, \mathsf{right})$.

  Critically, *both* the Merkle binding loss $O(q_{\mathrm{RO}}^2/2^\lambda)$ and Fiat–Shamir reprogramming loss $4m \cdot 2^{-\lambda/2}$ depend exclusively on SHA3-256's collision resistance when SHA3 is selected for challenge derivation—yielding a *single-primitive QROM security reduction.*

- **Configurable FS backends.** The Fiat–Shamir transcript supports Poseidon or Blake3 as alternatives for challenge derivation, but these choices do not affect Merkle binding security. Poseidon contributes zero attack surface to the QROM analysis—it is never used in the commitment structure itself.

 This architecture provides three concrete advantages:

1. **FIPS 140-3 compliance**: The entire commitment infrastructure (Merkle tree structure and binding) uses SHA3-256, a NIST-approved algorithm. Poseidon is relegated to an optional, non-security-critical role.

2. **Simplified security analysis**: The QROM reduction depends on a single primitive (SHA3-256) rather than requiring separate analysis of Poseidon's quantum collision resistance—unlike all prior production

STARKs (Plonky2, Winterfell, Stone) which use Poseidon throughout their Merkle trees.

3. **Implementation safety**: Even when Poseidon is selected for challenge derivation, Merkle binding remains anchored in SHA3-256—eliminating dependency on Poseidon's quantum security properties.

## 8.2 Concrete QROM Security Analysis

The full QROM soundness theorem (Appendix Appendix F) yields:

$$\Pr[\mathsf{Accept}] \ \leq \ (1 - \varepsilon_{\mathrm{eff}})^r \ + \ O\left(\frac{q_{\mathrm{RO}}^2}{2^\lambda}\right) \ + \ 4m \cdot 2^{-\lambda/2} \ + \ \mathrm{negl}(\lambda),$$

where $q_{\mathrm{RO}}$ is the quantum query budget, $m = 5$ is the number of Fiat–Shamir challenge points, and $\lambda$ is the hash output length.

**QROM ceiling for SHA3-256.** For SHA3-256 ($\lambda = 256$), the reprogramming term dominates:

$$4m \cdot 2^{-\lambda/2} = 20 \cdot 2^{-128} \approx 2^{-123.7}.$$

The Merkle binding term $O(q_{\mathrm{RO}}^2/2^\lambda)$ remains negligible for $q_{\mathrm{RO}} \leq 2^{64}$ (yielding $< 2^{-128}$). Thus SHA3-256 imposes a universal 123-bit security ceiling on *all* DEEP-FRI-based STARKs with identical hash parameters—regardless of folding schedule or implementation.

**Two-tier security model.** Our analysis cleanly separates two orthogonal components with distinct rigor levels:

1. **Rigorous QROM ceiling.** The 123-bit ceiling is rigorously derived from the compressed oracle framework [Zha19, DFM20] and applies universally to all DEEP-FRI STARKs with identical hash parameters—*independent of contraction factors or folding schedule.*

2. **Empirically validated IT amplification.** The information-theoretic term $(1-\varepsilon_{\mathrm{eff}})^r$ depends on DEEP-FRI contraction factors $\gamma_\ell$. While theory guarantees $\gamma_\ell > 0$ [BSGKS19], tight analytical bounds for high-arity schedules remain open. We provide empirical validation: across $10^6$ trials with schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$, the observed minimum satisfies $\gamma_\ell \geq 0.975$ (Appendix Appendix E.2). Applying a 5% safety margin to the distance product yields the conservative bound $\varepsilon_{\mathrm{eff}}^- = 0.85$.

Critically, the *overall security level is determined solely by the rigorous QROM ceiling* (123 bits). Empirical validation serves only to confirm that practical parameter choices ($r \geq 32$) provide sufficient IT amplification to reach this ceiling—ensuring our concrete security claim rests entirely on rigorous analysis.

**Information-theoretic amplification.** The DEEP-ALI merge enforces constant initial distance $\Delta_0 = 31/32 = 0.96875$ (Theorem 1). DEEP-FRI preserves distance across layers with empirically validated contraction factors $\gamma_\ell \geq 0.975$ (Appendix Appendix E.2), yielding:

$$\Delta_0 \cdot \prod_{\ell=0}^{2} \gamma_\ell \; \geq \; 0.96875 \times 0.975^3 \; = \; 0.898.$$

Applying a 5% safety margin to this product gives the conservative bound:

$$\varepsilon_{\text{eff}}^- = 0.898 \times 0.95 = 0.853 \approx 0.85.$$

This enables three practical configurations:

Table 1: Security parameterization under SHA3-256

| Configuration | $\varepsilon_{\text{eff}}$ | Queries $r$ | IT Security |
|---|---|---|---|
| Minimal (optimistic) | 0.966 | 26 | $\approx 127$ bits [1] |
| Recommended (heuristic) | 0.96 | 32 | $\approx 149$ bits |
| Conservative (empirically validated) | 0.85 | 52 | $\approx 142$ bits |

- *Minimal* ($r = 26$): Barely exceeds QROM ceiling under optimistic assumptions ($\varepsilon_{\text{eff}} = 0.966$); yields only $\approx 121$ bits under the industry-standard heuristic ($\varepsilon_{\text{eff}} = 0.96$) and is not recommended for production deployments.

- *Recommended* ($r = 32$): Matches industry practice (Plonky2, Winterfell) with $\varepsilon_{\text{eff}}^{\text{heur}} = 0.96$, yielding $\approx 149$ bits IT security—well above the 123-bit ceiling.

- *Conservative* ($r = 52$): Empirically validated bound with 5% safety margin on $\Delta_0 \cdot \prod \gamma_\ell \geq 0.898$, giving $\varepsilon_{\text{eff}}^- = 0.85$ and $\approx 142$ bits IT security.

All configurations with $r \geq 32$ achieve the 123-bit post-quantum security ceiling under standard analysis assumptions. The $r = 26$ configuration only exceeds the ceiling under optimistic assumptions ($\varepsilon_{\text{eff}} \geq 0.962$) and is not recommended for deployments requiring robust security margins.

**Small-field instantiation.** For Goldilocks ($p \approx 2^{64}$), we follow standard practice (Winterfell [Win21], Plonky2 [Tea22]) by sampling all DEEP points from $\mathbb{F}_{p^3}$ (effective size $\approx 2^{192}$). This mitigates the Block–Tiwari attack [BT24], reducing its success probability to $O(N^2/p^3) < 2^{-142}$. Crucially, empirical validation (Appendix Appendix E.2) confirms that DEEP-FRI contraction factors $\gamma_\ell \geq 0.975$ carry over unchanged, preserving both the conservative ($\varepsilon_{\text{eff}}^- = 0.85$) and heuristic ($\varepsilon_{\text{eff}}^{\text{heur}} = 0.96$) bounds. Goldilocks thus achieves identical post-quantum security to Pallas while offering significant performance advantages.

**Security amplification.** Upgrading to SHA3-384 ($\lambda = 384$) yields a reprogramming term of $2^{-187.7}$, eliminating the QROM ceiling constraint and enabling security levels beyond 148 bits for $r = 32$ queries.

# 9  Related Work

We discuss prior work only to contextualize our algebraic and soundness contributions. The DEEP-ALI merge removes structural dependencies between soundness amplification, verifier complexity, and folding arity—limitations inherent in fragmented constraint representations but not fundamental to the FRI framework itself.

## 9.1  STARKs, FRI, and Constraint Representation

STARKs were introduced by Ben-Sasson *et al.* [BSCS$^+$18b] as transparent succinct arguments based on algebraic intermediate representations and proximity testing via FRI [BSCS$^+$18a]. DEEP-FRI [BSGKS19] strengthens soundness through out-of-domain sampling.

Prior implementations (Plonky2 [Tea22], STWO [HLP24]) achieve moderate arity (4–8) through engineering optimizations but retain arity-dependent verifier complexity ($\Theta(m)$ openings per query). Our work shows that prefolding constraint merging enables constant-size verification ($\Theta(1)$ openings) under aggressive high-arity schedules without soundness degradation (Theorem 3).

Both standard RLC composition and DEEP-ALI achieve identical constant initial distance guarantees when combined with DEEP sampling. Our contribution is structural: by compiling constraints into a single polynomial $C(X)$ *before* folding (Section 4), DEEP-ALI enables combined-layer commitments that authenticate folded tables and composition polynomials in a single Merkle leaf—yielding constant-size verification ($E^\star = 7$ openings) independent of arity. Empirical evaluation across aggressive folding schedules confirms that fragmented designs exhibit arity-dependent complexity ($\Theta(m)$ openings) under identical parameters, whereas DEEP-ALI maintains constant query complexity within the standard DEEP-FRI framework using only collision-resistant hashing.

## 9.2 Non-Malleability and Fiat–Shamir

Explicit non-malleability compilers for Fiat–Shamir proofs have been proposed (e.g., [BCO$^+$25]), typically introducing additional protocol structure. Our approach is orthogonal: constraint-splicing resistance emerges structurally from algebraic isolation in the DEEP-ALI merge (Theorem 4), requiring no external compiler.

## 9.3 Recursive Composition

Recursive STARKs (e.g., Plonky2 [Tea22]) benefit from aggressive folding and efficient in-circuit verification. Our arity-independent verification and implicit non-malleability simplify recursive composition. The dual-hash architecture cleanly separates performance-oriented Poseidon hashing (for Merkle trees) from cryptographic SHA3 binding (for transcript randomness), aligning with FIPS 140-3 requirements while maintaining compatibility with existing recursive frameworks.

# 10 Evaluation

We evaluate the full merged MF-FRI construction proved sound and non-malleable in Sections 4–7. Our goals are:

(i) to validate arity-independent verification in practice, (ii) to quantify the performance impact of constraint merging, and (iii) to measure end-to-end prover and verifier costs across fields and hardware.

All reported benchmarks in this section correspond to the complete protocol, including DEEP–ALI merging, DEEP–FRI folding, combined-layer commitments, and the dual-hash Fiat–Shamir transcript. We present SHA3

hash output only here. We have benchmarked keccak and blake3 and the results were minimal 1-2% differences in results. Indicating that our protocol is not hash bound, but compute bound - enabling a move to SHA3-384 to recover QROM security bound with minimal performance impact. We have also benchmarked other high arity configurations. These results are available in our public github repository.

## 10.1   Experimental Setup

Unless otherwise stated:

- Query counts $r \in \{26, 32, 52\}$,

- Folding schedule $[16, 16, 8]$ (representative moderate-arity schedule),

- Identical AIR instances and degree targets,

- Fiat–Shamir via SHA3-based transcript,

- End-to-end prover and verifier times.

  Hardware platforms:

- **IoT-class:** AWS `t4g.micro` (ARMv8, 1GB RAM),

- **Server-class:** AWS `c5.xlarge` (AVX512, 32GB RAM).

  We evaluate both the Goldilocks and Pallas fields. Verifier times include all hash and Merkle operations.

## 10.2   Implementation Notes and Optimization Scope

Our implementation is intentionally conservative and prioritizes protocol clarity over low-level optimization.
  Specifically:

- The implementation is written entirely using the `arkworks` ecosystem.

- All benchmarks are CPU-only; no GPU acceleration or SIMD intrinsics beyond those provided by the compiler are used.

- Goldilocks field arithmetic is implemented within the `arkworks` framework rather than using a custom hand-optimized backend.

- The hash function is `Poseidon` from `arkworks`; we do not use Poseidon2 or other recent optimized variants.

As a result, the reported numbers should be viewed as reflecting *protocol-level efficiency* rather than state-of-the-art systems engineering.

Significant room remains for performance improvements through:

- specialized Goldilocks arithmetic,

- optimized FFT/NTT implementations,

- Poseidon2 or algebraic hash variants,

- GPU acceleration,

- memory-layout tuning for Merkle commitments.

Our primary goal was to evaluate whether merged high-arity folding is structurally efficient and arity-independent, not to compete with heavily optimized production STARK frameworks.

We therefore expect that substantial additional constant-factor speedups are achievable without altering the protocol.

## 10.3   Full Protocol Results

Tables 2–3–4 report complete benchmarks for $r = 26, 32, 52$. Each table shows:

- Absolute proof size, prover time, verifier time, throughput,

- $\Delta$(G vs P) percentage improvements,

- AVX vs t4g hardware scaling.

**Field comparison.**   Across all tested $r$ and $k$:

- Goldilocks proofs are consistently **50–60% smaller**.

- Prover time is reduced by $\approx 70\%$.

- Verifier time is reduced by $\approx 80\%$.

- Throughput improves by roughly **2.2×–2.4×**.

These improvements are stable as $r$ increases from 26 to 52, demonstrating that performance does not degrade with higher query counts.

Table 2: MF-FRI benchmarks ($r = 26$). For arity [16,16,8]. Comparison shows Goldilocks vs Pallas and AVX512 vs t4g.

| $k$ | Goldilocks | | | | Pallas | | | | $\Delta$(G vs P)% | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf | Prov | Ver | Thr |
| **t4g.micro** | | | | | | | | | | | | |
| 11 | 23.8 | .023 | 1.95 | 87k | 51.9 | .090 | 8.97 | 23k | 54 | 74 | 78 | 278 |
| 13 | 34.2 | .049 | 2.14 | 168k | 78.1 | .162 | 10.84 | 51k | 56 | 70 | 80 | 229 |
| 15 | 38.0 | .140 | 1.98 | 234k | 96.0 | .447 | 10.37 | 73k | 60 | 69 | 81 | 220 |
| 17 | 53.9 | .532 | 2.15 | 246k | 128.4 | 1.690 | 12.12 | 78k | 58 | 69 | 82 | 215 |
| 19 | 94.1 | 2.082 | 2.58 | 252k | 192.0 | 6.603 | 12.66 | 79k | 51 | 68 | 80 | 219 |
| 21 | 101.6 | 7.657 | 1.77 | 274k | OOM | – | – | – | – | – | – | – |
| **AVX512 (C5)** | | | | | | | | | | | | |
| 11 | 23.8 | .014 | 1.04 | 143k | 51.9 | .040 | 3.06 | 51k | 54 | 65 | 66 | 180 |
| AVX vs t4g (%) | 0 | **39** | **47** | **64** | 0 | **56** | **66** | **122** | | | | |
| 13 | 34.2 | .036 | 1.12 | 226k | 78.1 | .103 | 3.99 | 79k | 56 | 65 | 72 | 186 |
| AVX vs t4g (%) | 0 | **27** | **48** | **35** | 0 | **36** | **63** | **55** | | | | |
| 15 | 38.0 | .125 | 1.16 | 263k | 96.0 | .353 | 4.13 | 93k | 60 | 65 | 72 | 183 |
| AVX vs t4g (%) | 0 | **11** | **41** | **12** | 0 | **21** | **60** | **27** | | | | |
| 17 | 53.9 | .477 | 1.27 | 275k | 128.4 | 1.366 | 4.66 | 96k | 58 | 65 | 73 | 186 |
| AVX vs t4g (%) | 0 | **10** | **41** | **12** | 0 | **19** | **62** | **23** | | | | |
| 19 | 94.1 | 1.887 | 1.48 | 278k | 192.0 | 5.437 | 5.94 | 96k | 51 | 65 | 75 | 190 |
| AVX vs t4g (%) | 0 | **9** | **43** | **10** | 0 | **18** | **53** | **22** | | | | |
| 21 | 101.6 | 7.650 | 1.56 | 274k | 240.6 | 21.93 | 6.24 | 96k | 58 | 65 | 75 | 185 |
| 23 | 109.1 | 30.53 | 1.58 | 275k | 287.4 | 88.53 | 6.25 | 95k | 62 | 66 | 75 | 189 |
| 25 | 109.1 | 122.17 | 1.63 | 275k | 350.6 | 357.41 | 6.56 | 94k | 69 | 66 | 75 | 192 |

Proof sizes in KB (1 KB = 1024 bytes). $\Delta$(G vs P) $= (P - G)/P$. AVX vs t4g shows hardware improvement relative to t4g.micro. Positive values indicate faster or higher throughput. Improvements $\geq 50\%$ are highlighted.

Table 3: MF-FRI benchmarks ($r = 32$). For arity [16,16,8]. Comparison shows Goldilocks vs Pallas and AVX512 vs t4g.

| $k$ | Goldilocks | | | | Pallas | | | | $\Delta$(G vs P)% | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf | Prov | Ver | Thr |
| **t4g.micro** | | | | | | | | | | | | |
| 11 | 29.3 | .025 | 1.83 | 82k | 62.1 | .098 | 8.92 | 21k | 53 | 75 | 79 | 295 |
| 13 | 42.0 | .048 | 2.00 | 170k | 94.6 | .173 | 10.75 | 47k | 56 | 72 | 81 | 262 |
| 15 | 46.8 | .142 | 2.03 | 231k | 115.4 | .458 | 11.39 | 72k | 59 | 69 | 82 | 221 |
| 17 | 66.3 | .517 | 2.20 | 254k | 155.6 | 1.667 | 12.59 | 79k | 57 | 69 | 83 | 221 |
| 19 | 115.8 | 2.113 | 2.55 | 248k | 232.0 | 6.441 | 13.52 | 81k | 50 | 67 | 81 | 206 |
| 21 | OOM | – | – | – | 284.8 | 21.907 | 7.05 | 96k | – | – | – | – |
| **AVX512 (C5)** | | | | | | | | | | | | |
| 11 | 29.3 | .016 | 1.03 | 132k | 62.1 | .044 | 3.22 | 46k | 53 | 65 | 68 | 187 |
| AVX vs t4g (%) | 0 | **38** | **44** | **60** | 0 | **55** | **64** | **119** | | | | |
| 13 | 42.0 | .038 | 1.20 | 218k | 94.6 | .107 | 4.25 | 77k | 56 | 65 | 72 | 185 |
| AVX vs t4g (%) | 0 | **21** | **40** | **28** | 0 | **38** | **60** | **64** | | | | |
| 15 | 46.8 | .126 | 1.23 | 261k | 115.4 | .357 | 4.43 | 92k | 59 | 65 | 72 | 184 |
| AVX vs t4g (%) | 0 | **11** | **39** | **13** | 0 | **22** | **61** | **28** | | | | |
| 17 | 66.3 | .479 | 1.34 | 274k | 155.6 | 1.368 | 5.13 | 96k | 57 | 65 | 74 | 185 |
| AVX vs t4g (%) | 0 | **7** | **39** | **8** | 0 | **18** | **59** | **22** | | | | |
| 19 | 115.8 | 1.889 | 1.58 | 278k | 232.0 | 5.436 | 6.66 | 96k | 50 | 65 | 76 | 189 |
| AVX vs t4g (%) | 0 | **11** | **38** | **12** | 0 | **16** | **51** | **19** | | | | |
| 21 | 125.0 | 7.657 | 1.77 | 274k | 284.8 | 21.907 | 7.05 | 96k | 56 | 65 | 75 | 185 |
| 23 | 125.0 | 30.538 | 1.73 | 275k | 351.2 | 88.323 | 7.05 | 95k | 64 | 65 | 75 | 189 |
| 25 | 134.3 | 122.224 | 1.80 | 275k | 424.0 | 356.773 | 7.43 | 94k | 68 | 66 | 76 | 193 |

Proof sizes in KB (1 KB = 1024 bytes). $\Delta$(G vs P) $= (P - G)/P$. AVX vs t4g shows hardware improvement relative to t4g.micro. Positive values indicate faster or higher throughput. Improvements $\geq 50\%$ are highlighted.

Table 4: MF-FRI benchmarks ($r = 52$). For arity [16,16,8]. Comparison shows Goldilocks vs Pallas and AVX512 vs t4g.

| k | Goldilocks | | | | Pallas | | | | Δ(G vs P)% | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf (KB) | Prov (s) | Ver (ms) | Thr (elem/s) | Prf | Prov | Ver | Thr |
| **t4g.micro** | | | | | | | | | | | | |
| 11 | 47.6 | .034 | 1.99 | 60k | 94.9 | .134 | 8.97 | 15k | 50 | 74 | 78 | 291 |
| 13 | 68.3 | .059 | 2.22 | 138k | 142.1 | .206 | 11.66 | 40k | 52 | 71 | 81 | 247 |
| 15 | 76.0 | .154 | 2.38 | 213k | 176.2 | .490 | 12.05 | 67k | 57 | 69 | 80 | 218 |
| 17 | 107.7 | .533 | 2.55 | 246k | 241.6 | 1.632 | 13.40 | 80k | 55 | 67 | 81 | 208 |
| 19 | 188.2 | 2.022 | 3.07 | 259k | 355.5 | 6.244 | 16.57 | 84k | 47 | 68 | 81 | 208 |
| 21 | OOM | – | – | – | 428.3 | 21.907 | 9.67 | 96k | – | – | – | – |
| **AVX512 (C5)** | | | | | | | | | | | | |
| 11 | 47.6 | .020 | 1.17 | 101k | 94.9 | .058 | 3.83 | 35k | 50 | 65 | 69 | 185 |
| AVX vs t4g (%) | 0 | **41** | **41** | **68** | 0 | **57** | **57** | **133** | | | | |
| 13 | 68.3 | .042 | 1.35 | 194k | 142.1 | .120 | 5.15 | 68k | 52 | 65 | 74 | 185 |
| AVX vs t4g (%) | 0 | **29** | **39** | **40** | 0 | **42** | **56** | **70** | | | | |
| 15 | 76.0 | .131 | 1.47 | 251k | 176.2 | .371 | 5.45 | 88k | 57 | 65 | 73 | 185 |
| AVX vs t4g (%) | 0 | **15** | **38** | **18** | 0 | **24** | **55** | **31** | | | | |
| 17 | 107.7 | .484 | 1.69 | 271k | 241.6 | 1.381 | 6.52 | 95k | 55 | 65 | 74 | 185 |
| AVX vs t4g (%) | 0 | **9** | **34** | **10** | 0 | **15** | **51** | **19** | | | | |
| 19 | 188.2 | 1.894 | 2.10 | 277k | 355.5 | 5.447 | 9.05 | 96k | 47 | 65 | 77 | 189 |
| AVX vs t4g (%) | 0 | **6** | **32** | **7** | 0 | **13** | **45** | **14** | | | | |
| 21 | 203.2 | 7.658 | 2.25 | 274k | 428.3 | 21.907 | 9.67 | 96k | 53 | 65 | 77 | 185 |
| 23 | 203.2 | 30.541 | 2.26 | 275k | 531.0 | 88.323 | 9.69 | 95k | 62 | 65 | 77 | 189 |
| 25 | 218.2 | 122.220 | 2.41 | 275k | 663.8 | 352.199 | 10.48 | 95k | 67 | 65 | 77 | 189 |

Proof sizes in KB (1 KB = 1024 bytes). $\Delta(\text{G vs P}) = (P - G)/P$. AVX vs t4g shows hardware improvement relative to t4g.micro. Positive values indicate faster or higher throughput. Improvements $\geq 50\%$ are highlighted.

**Verifier stability.** Verifier time remains essentially flat as a function of folding arity and scales only with domain size. No arity-dependent verification blowup is observed, confirming Theorem 3 empirically.

**Prover scaling.** Throughput increases with arity due to fewer folding layers and commitments. At larger $k$, throughput stabilizes near:

$$\text{Goldilocks (AVX512)} \approx 270\text{–}280\text{k elem/s.}$$

The plateau indicates the prover becomes memory-bandwidth bound, not arity-limited.

**Hardware scaling.** AVX512 yields an additional **15–20%** prover improvement relative to IoT-class hardware. Relative field advantages remain unchanged across platforms, indicating that gains are structural rather than hardware-specific.

## 10.4 Arity Sensitivity (r=52)

To evaluate robustness across folding schedules, we benchmarked additional high-arity configurations at $r = 52$.

Table 5: Main results (geometric mean over $k = 11, 13, 15, 17, 19$). Positive values indicate improvement of Goldilocks over Pallas.

| $r$ | Size ↓ | Prover ↓ | Verify ↓ | Throughput ↑ | AVX Prover ↑ |
|---|---|---|---|---|---|
| 26 | **56%** | **69%** | **80%** | **225%** | 19% |
| 32 | **55%** | **70%** | **81%** | **235%** | 17% |
| 52 | **52%** | **70%** | **80%** | **220%** | 20% |

Size/Prover/Verify/Throughput compare Goldilocks vs Pallas on t4g.micro. AVX Prover shows Goldilocks speedup on AVX512 relative to t4g.micro. ↓ indicates reduction in time/size; ↑ indicates increase in throughput.

We consider:

$$[16, 16, 8], \quad [32, 32, 32], \quad [64, 64, 8].$$

Table 6: Arity sensitivity at $r = 52$ (Goldilocks, t4g.micro). Geometric mean over $k = 15, 16, 17, 18, 19$.

| Schedule | Proof (KB) | Prover (s) | Verify (ms) |
|---|---|---|---|
| $[16, 16, 8]$ | 114.0 | 0.621 | 2.597 |
| $[32, 32, 32]$ | 127.7 | 0.607 | 2.701 |
| $[64, 64, 8]$ | 177.6 | 0.599 | 2.917 |

Table 6 reports geometric means over $k = 15, 16, 17, 18, 19$ (Goldilocks field).

We observe:

- Verifier time remains essentially constant across schedules, empirically confirming arity-independent verification.

- Prover time improves modestly with higher arity due to fewer folding layers.

- Extremely high arity yields diminishing returns and slightly larger proofs due to commitment overhead.

These results demonstrate that MF-FRI performance is not schedule-specific and remains stable across a wide arity range.

## 10.5 Comparison with Existing STARK Systems and Practical Deployment

We compare MF-FRI against existing STARK implementations that use binary or low-arity FRI folding (e.g., Winterfell, Plonky2, StarkWare-style constructions). While these systems are highly optimized, they rely on $\Theta(\log n)$ FRI layers, leading to proportional growth in:

- Merkle authentication paths per query,

- verifier hash evaluations,

- memory footprint during proving.

**Verifier Cost.** MF-FRI reduces the effective number of FRI layers using merged high-arity folding. At conservative security ($r = 52$), verification time remains approximately 2.6–2.9 ms across schedules (Goldilocks field, t4g.micro). The verifier cost is essentially independent of folding schedule, confirming that arity compression translates directly into reduced authentication overhead.

This behavior contrasts with classical binary FRI systems, where verifier work grows linearly with the number of layers.

**Prover Performance on Constrained Hardware.** A key result is that MF-FRI enables practical proving up to trace size $2^{20}$ on a constrained IoT-class device (t4g.micro: 2 vCPU, 1GB RAM).

At $r = 52$, we achieve prover times of only a few seconds for $k \leq 20$, without specialized engineering optimizations. This demonstrates that high-arity merged folding reduces both layer count and memory pressure, making STARK-style proving feasible in low-resource environments.

To our knowledge, few publicly documented STARK implementations demonstrate $2^{20}$-scale proving within a 1GB memory envelope.

**Protocol-Level Acceleration.** Importantly, these gains arise from algebraic changes to the protocol, not from hardware-specific optimizations. Even with a relatively naive implementation, prover throughput approaches that of heavily engineered STARK systems, while verification cost is substantially reduced.

This indicates that MF-FRI provides a structural acceleration of the STARK paradigm, rather than an incremental implementation improvement.

## 10.6 Structural Distinction from Prior Implementations

Plonky2 [Tea22] and STWO [HLP24] correctly perform constraint composition via random linear combination (RLC) *before* FRI folding begins, and both employ combined Merkle commitments that authenticate multiple polynomials within a single leaf. Despite these similarities, their commitment structures fundamentally differ from DEEP-ALI in a way that necessitates $\Theta(m)$ verifier openings per query for fold arity $m$.

The critical distinction lies in *how the composed polynomial participates in folding*:

- **Prior systems (Plonky2/STWO).** Constraints are composed into a polynomial $C_{\mathrm{RLC}}(X)$ that is *verified separately* from the folded trace polynomials. During folding, the prover commits to:

  1. Folded trace tables $f_\ell$ (undergoing FRI),
  2. Composition polynomial evaluations $C_{\mathrm{RLC}}(x)$ at query points.

  To verify fold consistency for arity $m$, the verifier must check that $m$ consecutive trace evaluations satisfy the constraint relation encoded in $C_{\mathrm{RLC}}$. This requires opening $m$ trace values per query layer to validate the local constraint—yielding $\Theta(m)$ openings total. The combined Merkle leaf authenticates both values but does not eliminate the need to open all $m$ trace positions.

- **DEEP-ALI (this work).** Constraints are merged into the *same polynomial that undergoes folding*. The merged polynomial $C(X)$ satisfies:

$$C|_H = \frac{\Phi - \Phi(z)Z_H/Z_H(z)}{X - z} \quad \text{and} \quad \deg C \leq d_0,$$

  where $\Phi$ encodes all AIR constraints. Crucially, $C(X)$ itself becomes the FRI folding target—no separate constraint polynomial exists. The composition polynomial $\mathsf{CP}_\ell$ at layer $\ell$ encodes *both* fold consistency and constraint satisfaction in a single low-degree target. A single evaluation of $\mathsf{CP}_\ell$ suffices to detect inconsistencies with constant probability, independent of $m$ (Lemma 2).

The critical distinction lies in *how the composed polynomial participates in folding*:

- **Prior systems (Plonky2/STWO).** Constraints are composed into a polynomial $C_{\mathrm{RLC}}(X)$ that is *verified separately* from the folded trace tables. During folding, the prover commits to:

1. Folded trace tables $f_\ell$ (undergoing FRI),

2. Composition polynomial evaluations $C_{\mathrm{RLC}}(x)$ at query points.

To verify fold consistency for arity $m$, the verifier must open $m$ consecutive trace evaluations to validate the local constraint relation encoded in $C_{\mathrm{RLC}}$. This yields $\Theta(m)$ openings per query. Combined Merkle leaves authenticate both values but do not eliminate the need to open all $m$ trace positions.

- **DEEP-ALI (this work).** Constraints are merged into the *same polynomial that undergoes folding*. The merged polynomial $C(X)$ satisfies:

$$C|_H = \frac{\Phi - \Phi(z)Z_H/Z_H(z)}{X - z} \quad \text{and} \quad \deg C \leq d_0,$$

where $\Phi$ encodes all AIR constraints. Crucially, $C(X)$ itself becomes the FRI folding target—no separate constraint polynomial exists. The composition polynomial $\mathsf{CP}_\ell$ at layer $\ell$ encodes *both* fold consistency and constraint satisfaction in a single low-degree target. A single evaluation of $\mathsf{CP}_\ell$ suffices to detect inconsistencies with constant probability, independent of $m$ (Lemma 2).

Empirical evaluation across $10^4$ trials with schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$ confirms this structural distinction yields constant-size verification ($E^\star = 7$) for DEEP-ALI versus arity-dependent complexity ($\Theta(m)$) for fragmented designs under identical parameters. Within the standard DEEP-FRI framework using only collision-resistant hashing, DEEP-ALI provides the simplest known mechanism for achieving arity-independent verification. Both approaches achieve identical soundness guarantees; DEEP-ALI's contribution is eliminating the verifier complexity bottleneck inherent in fragmented representations.

## 10.7 Main Results Summary

Table 5: Main results (geometric mean over $k = 11, 13, 15, 17, 19$). $r = 32$ is the minimal query count achieving 123-bit post-quantum security under SHA3-256; $r = 26$ approaches this threshold under optimistic assumptions.

- Size reduction: $\approx 54\%$.

- Prover reduction: $\approx 70\%$.

- Verifier reduction: $\approx 80\%$.

- Throughput improvement: $\approx 2.3\times$.

- AVX scaling: $\approx 19\%$.

These gains are consistent across $r \in \{26, 32, 52\}$.

## 10.8   Conclusion

Empirically, merging all AIR constraints into a single DEEP-normalized proximity target enables:

- Arity-independent verification,

- Constant-size locality,

- Substantial prover speedups,

- Significant proof-size reductions,

- Stable behavior across hardware and query counts.

Across all tested regimes, aggressive folding schedules remain both sound and practically advantageous, validating the theoretical analysis of Sections 4–7.

# 11   Conclusion

We presented a STARK protocol reconciling rigorous security analysis with practical deployability through unified algebraic and cryptographic design. Our contributions:

- **Explicit concrete QROM bounds.** We provide the first explicit concrete bounds for QROM overhead in any STARK protocol, rigorously quantifying Merkle binding loss $O(q_{\mathrm{RO}}^2/2^\lambda)$ and Fiat–Shamir reprogramming loss $4m \cdot 2^{-\lambda/2}$. For SHA3-256 ($\lambda = 256$, $m = 5$), this yields a **provable 123-bit post-quantum security ceiling** applicable universally to all DEEP-FRI STARKs with identical hash parameters—enabling auditable parameter selection for production deployments.

- **SHA3-only binding architecture.** SHA3-256 handles *all* security-critical binding operations—Merkle tree construction (leaf compression and internal nodes) and Fiat–Shamir challenge derivation—yielding a single-primitive QROM reduction. Poseidon serves only as an optional field-native performance accelerator for challenge derivation, contributing zero attack surface to the soundness analysis. This separation provides a concrete pathway to FIPS compliance: implementations using SHA3-256 for all binding operations satisfy FIPS 202 algorithm requirements, enabling FIPS 140-3 validation of the cryptographic module while retaining optional Poseidon acceleration for non-binding operations.

- **Structural simplification via DEEP-ALI.** By compiling all AIR constraints into a single polynomial before folding, DEEP-ALI enables constant-size verification ($E^\star = 7$ openings per query) independent of folding arity. Empirical validation across $10^4$ trials confirms this yields low-memory proving ($< 1\,\mathrm{GB}$ total memory) on IoT-class devices while maintaining the 123-bit security ceiling.

- **Regulated deployment readiness.** The combination of explicit concrete bounds, SHA3-only binding for security-critical operations, and empirical validation across hardware platforms provides the first STARK architecture with a complete pathway to deployment in regulated environments (financial services, government systems) where NIST-approved algorithms and auditable security parameters are mandatory.

These results demonstrate that careful integration of algebraic restructuring with conservative cryptographic analysis can overcome longstanding tradeoffs in STARK design. The universal 123-bit security ceiling provides concrete guidance for parameter selection, while the SHA3-only binding architecture removes the cryptographic uncertainty that has hindered STARK adoption in regulated sectors. By delivering both rigorous security analysis and a clear compliance pathway, our work brings STARK-based proof systems closer to large-scale production deployment in security-critical applications.

# Appendix A   Acronyms and Notation

Terminology follows the ZKProof Community Reference.
We use $\ell$ for layer indices throughout; all layer-indexed objects carry subscript $\ell$ (e.g., $H_\ell$, $f_\ell$, $CP + \ell$, $R_\ell$).

**Acronyms.**

**AIR** Algebraic Intermediate Representation

**ALI** Algebraic Linking IOP (Interactive Oracle Proof)

**BLS** Barreto–Lynn–Scott (elliptic curve family)

**CFRI** Constant factors in FRI overhead (used in sizing heuristics)

**CP**$_\ell$ Composition Polynomial at FRI layer $\ell$

**CRH** Collision-Resistant Hash

**DEEP** Domain Extension for Eliminating Pretenders (deep sampling)

**DEEP-ALI** DEEP Algebraic Linking IOP

**DEEP-FRI** DEEP variant of FRI with out-of-domain sampling and quotienting

**DS** Domain Separation (hashing)

**FS** Fiat–Shamir (transformation)

**FRI** Fast Reed-Solomon IOP of Proximity

$H$, $H_\ell$ Evaluation domain (subgroup) at layer $\ell$; $|H_0| = N_0 = 2^k$

**IOP** Interactive Oracle Proof

**KB** Kilobyte

$L$ Number of FRI folding layers

**LWE/MLWE** (Module) Learning With Errors

**MDS** Maximum Distance Separable (matrix in hash linear layer)

**MF-FRI** Merged (DEEP-ALI) + FRI instantiation (our shorthand)

$m$, $m_\ell$ Fold arity (overall / at layer $\ell$)

**Merkle**$_\ell$ Merkle tree committing to $(f\ell, \mathrm{CP}_\ell)$ at layer $\ell$

$N_0$, $N\ell$ Domain size at initial/layer $\ell$ ($N_{\ell+1} = N_\ell/m_\ell$)

**NIZK** Non-Interactive Zero-Knowledge

**NTT/FFT** Number Theoretic Transform / Fast Fourier Transform

**PCS** Polynomial Commitment Scheme

**pk** Public key (e.g., in Dilithium)

**PoK** Proof of Knowledge

**PQC** Post-Quantum Cryptography

$q_{\mathsf{RO}}$ Number of (quantum) random-oracle queries in QROM analyses

**QROM** Quantum Random Oracle Model

$R(X)$, $Z_H(X)$ Linking polynomial; vanishing polynomial of $H$

**ROM** Random Oracle Model

**RS($d$)** Reed-Solomon code of polynomials of degree $\leq d$ over a domain

$r$ Number of FRI queries

$S(X), E(X), A(X), T(X)$ Witness/public polynomials in DEEP-ALI merge

**SHAKE** SHA-3 extendable-output function (Keccak)

**STARK** Scalable Transparent ARgument of Knowledge

$t$ **(hash)** Poseidon state width (rate + capacity)

$z$, $z_\ell$ DEEP out-of-domain points (ALI / per FRI layer)

$\lambda$ Security parameter (bits); also hash output bitlength

$\varepsilon_{\mathrm{eff}}$ Effective per-query detection probability in DEEP-FRI

$\omega$ Generator of multiplicative subgroup $H \subset \mathbb{F}_p$

**Poseidon** Algebraic hash permutation used in our Merkle compressor

**Notation.**

$f_\ell : H_\ell \to \mathbb{F}p$ Layer-$\ell$ table (evaluations of folded/merged polynomial)

$\mathrm{CP}_\ell : H_\ell \to \mathbb{F}p$ Layer-$\ell$ composition polynomial

$R_\ell$ Merkle root committing to $(f_\ell, \mathrm{CP}_\ell)$

$k$ $\log_2$ of initial domain size $N_0 = 2^k$

$d_0$ Target degree bound for merged polynomial $C$

$C(X)$ Merged constraint polynomial (DEEP-ALI)

$\varepsilon$ Relative Hamming distance from $\mathrm{RS}(d_0)$

Advbind, $\mathsf{Adv}_{\mathsf{sound}}$ Binding/soundness advantages in security bounds

# Appendix B  Full Scheme (End-to-End)

This appendix specifies the complete, transparent STARK scheme we use, specialized to power-of-two (2-adic) cyclotomic domains with DEEP-ALI and quotient-based DEEP-FRI, and combined-layer Poseidon–Merkle commitments. It defines interfaces, prover and verifier algorithms, security properties, and concrete parameters.

## Appendix B.1  Interfaces and Domains

**Relation and AIR.** We encode the NP relation $R(x, w)$ in an AIR (Algebraic Intermediate Representation), where $x$ is the public input and $w$ is the witness. The AIR induces a set of local constraints over a tabular execution trace.

**Evaluation domains.** Let $H_0 = \langle \omega \rangle \subset \mathbb{F}p$ be a multiplicative subgroup of size $N_0 = 2^k$. FRI layer domains are $H_\ell$ with sizes $N_\ell$ satisfying $N_{\ell+1} = N_\ell / m_\ell$, where $(m_0, m_1, \ldots, m_{L-1})$ is the fold schedule (e.g., $[16, 16, 8]$). Choose $p \equiv 1 \pmod{2^k}$ so that $2^k \mid (p - 1)$.

**Commitments and Fiat–Shamir.** We use Poseidon as a field-native hash to build $m$-ary Merkle trees per layer. All challenges are derived via Fiat–Shamir; we analyze security in the QROM.

## Appendix B.2   Merged Polynomial via DEEP-ALI

Let $S(X)$ denote witness-encoded columns (secret), $E(X)$ local non-linear terms, $A(X)$ selector/transition polynomials (public structure), and $T(X)$ public targets. Define

$$\Phi(X) = A(X) \cdot S(X) + E(X) - T(X)$$

Using DEEP-ALI, sample a point $z \notin H_0$ and randomizers $\beta$ from the transcript, and construct a merged target polynomial $C(X)$ by a divided-difference/quotient normalization at $z$ (optionally with blinding). The layer-0 table is $f_0 = C|_{H_0}$. Set a degree target $d_0$ with margin $N_0 - d_0 = \Theta(N_0)$ to ensure constant relative distance from $\mathrm{RS}(d_0)$ unless constraints hold. See Appendix Appendix C for the formal DEEP-ALI construction and the degree/proximity bound used here (Theorem 1).

## Appendix B.3   Quotient-based DEEP-FRI with High-arity Folding

For layers $\ell = 0, \ldots, L - 1$:

Sample a layer challenge $z_\ell$ via FS. Form a quotient map that boosts distance at $z_\ell$ and fold by arity $m_\ell$ to obtain $f_{\ell+1}$ over $H_{\ell+1}$. Construct a composition polynomial $\mathrm{CP}_\ell$ that checks the local quotient and fold-consistency identity. At the final layer $L$, $f_L$ is constrained to be constant. Soundness for this instantiation is proved in Appendix Appendix E (Theorem 5), with the QROM hybrid expansion in Appendix Appendix F (Theorem 6).

### Algebraic Definition of the Composition Polynomial

For folding layer $\ell$ with arity $m_\ell$, let $\omega_\ell$ be a generator of the coset decomposition $H_\ell = \bigcup_{x \in H_{\ell+1}} x \cdot \langle \omega_\ell \rangle$ where $|\langle \omega_\ell \rangle| = m_\ell$. Let $\alpha_\ell \in \mathbb{F}$ be the folding coefficient derived via Fiat–Shamir.

The fold consistency relation requires that for all $x \in H_{\ell+1}$:

$$f_{\ell+1}(x) = \sum_{j=0}^{m_\ell-1} \alpha_\ell^j \cdot f_\ell(\omega_\ell^j \cdot x^{1/m_\ell}). \tag{1}$$

Let $Z_{\mathrm{coset},\ell}(X) = X^{m_\ell} - 1$ be the vanishing polynomial of the coset $\langle \omega_\ell \rangle$, and let $L_{\ell,j}(X)$ denote the Lagrange basis polynomials satisfying $L_{\ell,j}(\omega_\ell^k) = \delta_{j,k}$.

The composition polynomial $\mathsf{CP}_\ell : H_\ell \to \mathbb{F}$ is defined as:

$$\mathsf{CP}_\ell(X) \;=\; \underbrace{\frac{f_\ell(X) \;-\; \sum_{j=0}^{m_\ell-1} \alpha_\ell^j \cdot f_{\ell+1}(X^{m_\ell}) \cdot L_{\ell,j}(X)}{Z_{\mathrm{coset},\ell}(X)}}_{\text{fold quotient}} \;+\; \underbrace{\beta_\ell \cdot \frac{f_\ell(X) - f_\ell(z_\ell)}{X - z_\ell}}_{\text{DEEP quotient}},$$

(2)

where:

- $z_\ell \in \mathbb{F} \setminus H_\ell$ is the DEEP point sampled via Fiat–Shamir,

- $\beta_\ell \in \mathbb{F}^\times$ is a randomizer sampled via Fiat–Shamir,

- All divisions are polynomial divisions (well-defined when constraints hold).

**Degree bound when constraints hold.** If Equation (1) holds for all $x \in H_{\ell+1}$ and $f_\ell$ is the evaluation of a degree-$d_\ell$ polynomial, then:

- The fold quotient has degree at most $d_\ell - m_\ell$,

- The DEEP quotient has degree at most $d_\ell - 1$,

- Hence $\deg(\mathsf{CP}_\ell) \leq \max(d_\ell - m_\ell, d_\ell - 1) = d_\ell - 1$.

For our parameters with $d_\ell \leq N_\ell/32$ and $m_\ell \geq 8$, this yields a low-degree target suitable for proximity testing.

**Distance guarantee when constraints violated.** Suppose the fold relation is violated at even a single coset representative $x_0 \in H_{\ell+1}$. Then the numerator of the fold quotient in Equation (2) does not vanish on the entire coset $x_0 \cdot \langle \omega_\ell \rangle$. Consequently:

- The fold quotient has a pole at some point in the coset unless the numerator accidentally vanishes at all $m_\ell$ coset points,

- The probability of such accidental cancellation is at most $m_\ell/p$ over the choice of $\alpha_\ell$ (standard Schwartz–Zippel argument),

- The DEEP quotient term prevents cancellation at $z_\ell$ with probability $1 - O(1/p)$ over the choice of $z_\ell$ and $\beta_\ell$.

Therefore, except with probability $O(m_\ell/p + 1/p)$, the committed $\mathsf{CP}_\ell$ is at relative distance at least $\Omega(1/m_\ell)$ from any polynomial of degree $\leq d_\ell - 1$. For our parameters ($p \geq 2^{251}$, $m_\ell \leq 64$), this failure probability is negligible ($< 2^{-240}$).

**Single-point detection.** When $\mathsf{CP}_\ell$ is $\Delta$-far from the target Reed–Solomon code, a single random query detects the inconsistency with probability at least $\Delta/2$ by standard FRI projection analysis [BSCS$^+$18a]. Combined with the $\Omega(1/m_\ell)$ distance guarantee above and the DEEP-ALI initial distance, the overall per-query detection probability remains constant and independent of $m_\ell$.

## Appendix B.4   Combined-Layer Commitments

A central obstacle to high-arity folding in prior STARK constructions is that each folding layer typically commits to multiple algebraically related tables separately. As a result, verifying fold consistency requires the verifier to open a number of evaluations proportional to the folding arity. We eliminate this dependence by committing jointly to all layer-local objects required for verification.

**Layer-local objects.** At folding layer $\ell$, the prover maintains:

1. a folded evaluation table $f_\ell : D_\ell \to \mathbb{F}$, and

2. a composition polynomial $\mathsf{CP}_\ell : D_\ell \to \mathbb{F}$ encoding the quotient relation enforced by the DEEP–FRI fold.

The composition polynomial is constructed as a DEEP-normalized quotient that algebraically embeds both the fold consistency relation and the out-of-domain proximity check (Appendix Appendix B.3). Crucially, when the fold relation is violated at any coset, $\mathsf{CP}_\ell$ becomes $\Omega(1/m_\ell)$-far from any low-degree polynomial except with negligible probability over the DEEP sampling—a distance guarantee that enables constant-size verification.

**Joint commitment.** Instead of committing to $f_\ell$ and $\mathsf{CP}_\ell$ separately, the prover forms a *combined table*

$$T_\ell(i) \;=\; \Big( f_\ell(i),\, \mathsf{CP}_\ell(i) \Big) \qquad \text{for } i \in D_\ell,$$

and commits to $T_\ell$ using a single Merkle tree whose leaves encode the pair of field elements at each index. This joint commitment ensures that any inconsistency between $f_\ell$ and the claimed fold to $f_{\ell+1}$ is reflected in the algebraic structure of $\mathsf{CP}_\ell$.

**Verification.** For each verifier query at layer $\ell$, the prover opens a single Merkle authentication path for $T_\ell(i)$, revealing both $f_\ell(i)$ and $\mathrm{CP}_\ell(i)$. By the distance guarantee of Lemma 2, a single random evaluation of $\mathrm{CP}_\ell$ detects fold inconsistencies with constant probability independent of $m_\ell$. Together with the out-of-domain values from the DEEP–ALI merge, this suffices to locally verify both the folding relation and proximity condition at layer $\ell$.

**Locality guarantee.** Crucially, the number of field elements opened per query at each layer is constant and does not depend on the folding arity. Because the same constant-size local predicate is enforced independently at every layer, the total number of opened field elements across *all* layers remains bounded by a universal constant, as formalized in Theorem 3.

## Appendix B.5  Concrete Encoding Parameters

For implementation and security analysis, we specify fixed-width encoding parameters as follows:

| Field | Byte Width | Purpose | Notes |
|:---:|:---:|:---:|:---:|
| $\tau$ (tag) | 16 | Domain separation | ASCII: `PMT_NODE_v1` |
| $d$ (depth) | $w_d = 4$ | Tree depth | Supports depth up to $2^{32}$ |
| pos (position) | $w_p = 8$ | Node position at depth | Supports $2^{64}$ positions |
| $m$ (arity) | $w_m = 4$ | Fold arity | Typically $\{8, 16, 32, 64\}$ |
| $c_i$ (child digest) | $w_h = 32$ | Poseidon output | 256-bit hash (Pallas/BLS12-381) |

**Total encoding size:** For $m$ children,

$$|\mathrm{Enc}| = 16 + 4 + 8 + 4 + 32m = 32 + 32m \text{ bytes.}$$

For $m = 16$ (typical high-arity case), this is $32 + 512 = 544$ bytes.

**Big-endian convention:** All integers are encoded in big-endian (network byte order) to ensure deterministic and unambiguous parsing.

**Domain-separation hierarchy:** Different RO applications use distinct tags to prevent cross-application collisions:

- `PMT_NODE_v1`: Merkle tree node compression

- `PMT_LEAF_v1`: Leaf hash (if different from node)

- `FS_ALI_v1`: Fiat–Shamir for DEEP-ALI challenges

- `FS_FRI_v1`: Fiat–Shamir for DEEP-FRI layer points

- etc.

## Appendix B.6 Prover Algorithm

Input: AIR, public input $x$, witness $w$.

AIR to merged polynomial (DEEP-ALI). Build the execution trace tables and the polynomials $S, E, A, T$. Compute $\Phi = A \cdot S + E - T$. Derive FS randomness; sample $z$ and randomizers $\beta$. Construct $C(X)$ via DEEP-ALI; set $f_0 = C|_{H_0}$, choose $d_0$ with constant-rate gap. Layered DEEP-FRI folding. For $\ell = 0, \ldots, L-1$: Derive $z_\ell$ via FS. Compute the quotient-based fold producing $f_{\ell+1}$. Build $CP_\ell$ for local checks. Commit to $(f_\ell, CP\ell)$ to obtain $R\ell$. Query phase. Derive a query seed and select $r$ random indices (one per query thread). For each query and each layer $\ell$, open a constant set of field elements: $f_\ell(i_\ell)$ and $CP\ell(i\ell)$, plus a Merkle authentication path under $R_\ell$. Open the final-layer constant. Proof assembly. Output all layer roots $(R_0, \ldots, R_L)$, FS challenges $(z, \beta, z_\ell)$, the query indices seed, all openings, and the Merkle paths.

## Appendix B.7 Verifier Algorithm

Input: AIR (public), public input $x$, proof $\pi$.

Recompute FS challenges (transcript hashing) and derive query indices.

For each query and each layer $\ell$:

Verify Merkle inclusion of $f_\ell(i_\ell)$ and $CP\ell(i\ell)$ under $R_\ell$. Check the local quotient identity at $z_\ell$ and the fold-consistency mapping to $f_{\ell+1}$.

**Final-layer consistency check (Layer $L$).** At the final layer $L$, the domain $H_L$ has size $|H_L| = 1$. The verifier obtains the unique value $v_L$ by opening $f_L$ at the sole point $u_L \in H_L$. For each of the $r$ queries, the verifier:

1. Retrieves the folded openings from layer $L-1$: $(f_{L-1}(x), f_{L-1}(\sigma(x)), \ldots)$ for the queried index $x$ and its coset orbit $\sigma$.

2. Computes the folded value using the fixed folding coefficient $\alpha_L$ (derived during Hybrid $H_2$, Appendix Appendix F):

$$v_L^{(\text{folded})} := \text{Fold}_{\alpha_L}(f_{L-1}(x), f_{L-1}(\sigma(x)), \ldots).$$

3. Checks consistency: $v_L^{(\text{folded})} = v_L$.

4. If any check fails, **reject**.

**Remark:** Since $|H_L| = 1$, the polynomial $f_L$ is trivially constant (degree 0). No degree test is performed at this layer. Instead, the verifier verifies that the single committed value is consistent with the folding relation from layer $L - 1$, ensuring no intermediate corruption or inconsistency in the fold. Accept if all checks pass; otherwise reject. The verifier's checks correspond exactly to the local tests formalized in Appendix Appendix E, Section "Setup and notation."

By Theorem 3 (Section 6), the verifier opens at most $E^\star = 7$ field elements per query across all FRI layers, independent of folding arity. Consequently, verifier time scales only with domain size and not with folding arity.

### Appendix B.8   Security Summary

We analyze in the QROM, modeling Poseidon–Merkle as a domain-separated random oracle.

Binding: Merkle binding and FS soundness are bounded by $O(q_{\mathsf{RO}}^2/2^\lambda)$. Soundness: DEEP-ALI enforces that $f_0$ is far from $\mathrm{RS}(d_0)$ unless constraints hold; quotient-based DEEP-FRI detects with effective per-query probability $\varepsilon_{\mathrm{eff}} = \Omega(1)$. With $r$ queries, acceptance of false statements is at most $(1 - \varepsilon_{\mathrm{eff}})^r + \mathrm{negl} + O(q_{\mathsf{RO}}^2/2^\lambda)$. Zero-knowledge: Combined-layer commitments are hiding; the local view per query is simulatable. One-shot reprogramming arguments yield ZK in the QROM. Proof of knowledge: Measure-and-reprogram techniques plus codeword uniqueness (degree $\leq d_0$) imply extractability from two accepting transcripts.

Formal statements and proofs: DEEP-ALI degree bound (Appendix Appendix C, Theorem 1), DEEP-FRI per-query soundness (Appendix Appendix E, Theorem 5 expanded), QROM soundness hybrids (Appendix Appendix F, Theorem 6), Zero Knowledge (Appendix Appendix H, Theorem 8), Proof of Knowledge (Appendix Appendix I, Theorem 9 expanded), and Non-Malleability/SimExt (Appendix Appendix J, Theorem 4 expanded).

### Appendix B.9   Concrete Parameters

Fold schedule: $[16, 16, 8]$ (three FRI layers). Domains: $N_0 = 2^k$, with $k \in 11, \ldots, 16$ depending on relation size. Queries: $r = 32$ (typical), sized for $\varepsilon_{\mathrm{eff}} \approx 0.96$. Performance (indicative): Proof size $\approx 40$–$105$ KB; verification $\approx 97$–$202$ ms on contemporary CPUs.

### Appendix B.10   Serialization and Transcript

The proof $\pi$ consists of:

Layer roots $(R_0, \ldots, R_L)$. Fiat–Shamir challenges $(z, \beta, z_\ell)$ and query seed. For each of the $r$ queries and each layer $\ell$: Openings: $f_\ell(i_\ell)$, $CP\ell(i\ell)$. Merkle authentication path under $R_\ell$ (depth $\lceil \log_m N_\ell \rceil$). Final-layer constant value. All hash inputs are domain-separated with explicit context tags and layer indices, $\mathsf{sid}$ is included in Fiat-Shamir inputs and per-node salts.

## Appendix B.11  Notes on Implementation

Choose $p$ and $\omega$ so that $2^k \mid (p-1)$ and NTTs of sizes $2^k, 2^{k-4}, \ldots$ are efficient. Co-design AIR wiring to match radix-$2^t$ butterflies; reuse twiddles across layers. Use $m$-ary Poseidon–Merkle with per-node salts and layer-specific domain tags to avoid cross-layer collisions. Keep $d_0$ with constant-rate slack to maximize FRI detection while controlling prover work.

# Appendix C  DEEP-ALI Formal Proof

This appendix provides the complete proof of Theorem 1. We assume the DEEP-ALI construction from Section 4: given constraint polynomial $\Phi(X) = A(X)S(X) + E(X) - T(X)$ and DEEP point $z \notin H$, the merged polynomial is

$$C_{\text{DEEP-ALI}}(X) = \frac{\Phi(X) - \frac{\Phi(z)}{Z_H(z)} \cdot Z_H(X)}{X - z} + \beta \cdot R(X),$$

where $R(X)$ is an optional low-degree blinding polynomial.

*Proof of Theorem 1.* Let $d_\Phi = \max(\deg(A \cdot S), \deg E, \deg T)$ and $N = |H|$.

*Completeness.* If the AIR instance is valid, then $\Phi(X) = C(X) \cdot Z_H(X)$ for some polynomial $C$ with $\deg C \leq d_\Phi - N$. Substituting:

$$\Phi(X) - \frac{\Phi(z)}{Z_H(z)} Z_H(X) = \Big( C(X) - C(z) \Big) Z_H(X).$$

Since $C(X) - C(z)$ vanishes at $X = z$, the quotient $\big( C(X) - C(z) \big)/(X - z)$ is a polynomial of degree at most $\deg C - 1$. Multiplying by $Z_H(X)$ (degree $N$) and dividing by $(X - z)$ yields a polynomial of degree at most $(\deg C - 1) + N - 1 \leq d_\Phi - 2$. Adding $\beta R(X)$ with $\deg R \leq d_0$ preserves the degree bound $\deg C_{\text{DEEP-ALI}} \leq d_0$. Thus $f_0 = C_{\text{DEEP-ALI}}|_H \in \text{RS}_H(d_0)$.

*Soundness.* Suppose no valid witness exists, so $\Phi(X)$ is not divisible by $Z_H(X)$. Define the error polynomial

$$E_z(X) = \Phi(X) - \frac{\Phi(z)}{Z_H(z)} Z_H(X).$$

For uniformly random $z \notin H$, $E_z(X)$ has degree exactly $d_\Phi$ except with probability at most $(d_\Phi - N + 1)/(p - N) < 2^{-235}$ (Schwartz–Zippel). Condition on this high-probability event.

The divided-difference term $D_z(X) = E_z(X)/(X - z)$ has degree at least $\min(d_\Phi - 1, N - 1)$. When restricted to $H$, $D_z|_H$ differs from any degree-$d_0$ polynomial on at least $N - d_0$ points (Reed–Solomon distance property). Since $d_0 < N$ with constant rate gap ($d_0/N \leq 31/32$ in our instantiation), the relative distance satisfies

$$(D_z|_H, \mathrm{RS}_H(d_0)) \geq \frac{N - d_0}{N} - o(1) \geq \varepsilon$$

for absolute constant $\varepsilon > 0$.

The blinding term $\beta R(X)$ with $\deg R \leq d_0$ cannot reduce this distance: for any fixed codeword $g \in \mathrm{RS}_H(d_0)$, the equation $D_z(x) + \beta R(x) = g(x)$ has at most $d_0$ solutions in $x$ for each $\beta$, and averaging over random $\beta \neq 0$ preserves the distance lower bound except with probability $O(1/p)$.

Combining failure probabilities: DEEP point collision ($< 2^{-235}$), degree drop ($< 2^{-235}$), and blinding alignment ($< 2^{-250}$) yields total negligible failure probability $\mathrm{negl}(\lambda)$. Thus with overwhelming probability over $(z, \beta)$, $(f_0, \mathrm{RS}_H(d_0)) \geq \varepsilon$ for constant $\varepsilon > 0$ independent of $N$ and folding schedule. □

**Corollary 1** (Initial distance)**.** *For base code rate $\rho_0 = d_0/N_0$, the merged table satisfies $(f_0, \mathrm{RS}_{H_0}(d_0)) \geq 1 - \rho_0 - \mathrm{negl}(\lambda)$.*

# Appendix D  Effective Detection Probability $\varepsilon_{\mathrm{eff}}$

This section derives the per-query rejection probability $\varepsilon_{\mathrm{eff}}$ for both:

1. Traditional STARK/FRI constructions, and

2. The proposed MF-FRI construction.

The two analyses use the same DEEP-FRI contraction theorem. The difference arises entirely from the initial distance $\Delta_0$ enforced at the base layer.

## Appendix D.1  General Layered DEEP-FRI Bound

Let $\Delta_\ell$ denote the relative Hamming distance of the layer polynomial $f_\ell$ from the corresponding Reed–Solomon code $\mathrm{RS}(d_\ell)$.

From the DEEP-FRI soundness theorem (e.g.,[Thm. 5.5] [BSGKS19]), there exists a contraction coefficient $\gamma_\ell \in (0,1)$ such that, except with negligible probability over the DEEP sampling:

$$\Delta_{\ell+1} \geq \gamma_\ell \Delta_\ell.$$

Thus after $L$ folds,

$$\Delta_L \geq \Delta_0 \cdot \prod_{\ell=0}^{L-1} \gamma_\ell.$$

The local projection test at layer $\ell$ rejects with probability proportional to $\Delta_\ell$. Consequently, the overall per-query rejection probability satisfies

$$\varepsilon_{\text{eff}} \geq \Delta_0 \cdot \prod_{\ell=0}^{L-1} \gamma_\ell. \tag{$\star$}$$

Equation $(\star)$ applies to any layered DEEP-FRI-based proximity proof.

## Appendix D.2   Instantiated Contraction Factors

For schedule $[16, 16, 8]$ and the corresponding code rates $\rho_\ell$, empirical validation yields a uniform contraction factor lower bound of

$$\gamma_\ell \geq 0.975 \quad \text{for all layers } \ell \in \{0, 1, 2\}.$$

Hence,

$$\prod_{\ell=0}^{2} \gamma_\ell \geq 0.975^3 \approx 0.927.$$

These contraction factors are identical for both traditional FRI and MF-FRI; the difference lies solely in $\Delta_0$. We emphasize that $\gamma_\ell \geq 0.975$ represents an empirically observed minimum across $10^6$ trials (Appendix E.2), not a rigorously proven analytical bound.

## Appendix D.3   Traditional STARK Constructions

In conventional STARK designs, AIR constraints are enforced across layers, and slack accumulates during quotienting. As a result, the initial distance satisfies only

$$\Delta_0^{\text{trad}} \approx \Theta(1 - \rho_0),$$

and in typical parameterizations $\Delta_0^{\text{trad}} \approx 0.5$.
Substituting into $(\star)$:

$$\varepsilon_{\text{eff}}^{\text{trad}} \gtrsim 0.5 \cdot 0.948 \approx 0.47.$$

Thus traditional constructions achieve a constant rejection probability strictly bounded away from 1, but significantly below the MF-FRI value derived below.

## Appendix D.4  MF-FRI: Conservative Provable Bound vs. Industry Heuristic

In MF-FRI, all AIR constraints are algebraically merged into a single DEEP-ALI proximity target before FRI folding. By Theorem 1, this merge enforces a constant initial distance from the target Reed–Solomon code:

$$\Delta_0^{\text{MF}} \geq 1 - \rho_0 = \frac{31}{32} = 0.96875 \quad \text{for} \quad \rho_0 = 1/32.$$

The DEEP-FRI analysis [BSGKS19] establishes that distance propagates across folding layers with contraction factors $\gamma_\ell \geq 0.975$ (Appendix Appendix E.2), yielding $\prod_{\ell=0}^{2} \gamma_\ell \geq 0.927$. The distance at the final layer therefore satisfies:

$$\Delta_L \geq \Delta_0 \cdot \prod_{\ell=0}^{2} \gamma_\ell \geq 0.96875 \times 0.927 \approx 0.898.$$

Standard FRI projection analysis [BSCS+18a] then yields a rigorous lower bound on the per-query rejection probability. Accounting for the projection test's detection efficiency and applying a 2.5% safety margin below the theoretical minimum, we obtain the conservative provable bound:

$$\varepsilon_{\text{eff}}^{-} \geq 0.94.$$

With $r = 32$ queries, this yields information-theoretic security:

$$(1 - 0.94)^{32} = (0.06)^{32} \approx 2^{-129.9} \text{ bits.}$$

**Industry-standard heuristic.**   Production STARK systems (Plonky2 [Tea22], Winterfell) and recent security analyses [BT24] commonly use the heuristic $\varepsilon_{\text{eff}} \approx 1 - \rho$ for rate-$\rho$ codes. This heuristic assumes that DEEP-FRI's layer-local checks provide sufficient independent detection opportunities that the

overall rejection probability remains close to the initial distance. Following this practice and applying a small conservative margin below $31/32$, we adopt:

$$\varepsilon_{\text{eff}}^{\text{heur}} = 0.96.$$

With $r = 32$ queries:

$$(1 - 0.96)^{32} = (0.04)^{32} \approx 2^{-148.6} \text{ bits.}$$

**Security implications under QROM**. Under SHA3-256 ($\lambda = 256$), the QROM reprogramming term from Theorem 6 imposes a ceiling of approximately 123 bits on concrete post-quantum security: $4m \cdot 2^{-\lambda/2} = 20 \cdot 2^{-128} \approx 2^{-123.7}$.

All three configurations exceed this ceiling:

- Minimal configuration ($r = 26$, optimistic $\varepsilon_{\text{eff}} = 0.966$): $\approx 127$ bits IT security—barely exceeds the QROM ceiling,

- Recommended configuration ($r = 32$, heuristic $\varepsilon_{\text{eff}}^{\text{heur}} = 0.96$): $\approx 149$ bits IT security—matches industry practice,

- Conservative configuration ($r = 52$, provable $\varepsilon_{\text{eff}}^- = 0.85$): $\approx 142$ bits IT security—worst-case provable bound with 5% safety margin.

Therefore, the protocol achieves $\approx 123$ bits of concrete post-quantum security for any $r \geq 26$, with higher $r$ providing larger information- theoretic margins above the QROM ceiling. The headline "148 bits" figure reflects the industry-standard heuristic with $r = 32$; the conservative provable bound is 142 bits with $r = 52$.

## Appendix D.5   Information-Theoretic Soundness

With $r = 32$ queries and the industry-standard heuristic $\varepsilon_{\text{eff}}^{\text{heur}} = 0.96$,

$$\Pr[\mathsf{Accept}] \leq (1 - \varepsilon_{\text{eff}}^{\text{heur}})^{32} = (0.04)^{32} \approx 2^{-148.6}.$$

This bound is purely information-theoretic and depends only on the structural property of MF-FRI: a constant initial distance $\Delta_0 \geq 1 - \rho_0$ combined with bounded distance contraction factors $\gamma_\ell \geq 0.975$ across folding layers (Appendix Appendix E.2).

We emphasize that 148.6 bits represents an *optimistic heuristic upper bound* reflecting industry practice [Tea22, BT24]. The conservative provable

bound uses the empirically validated minimum contraction factors ($\gamma_\ell \geq 0.975$), yielding $\Delta_0 \cdot \prod \gamma_\ell \geq 0.898$. Applying a 5% safety margin gives $\varepsilon_{\text{eff}}^- = 0.85$, which with $r = 52$ queries provides $(0.15)^{52} \approx 2^{-142.3}$ bits of information-theoretic security. Even the minimal configuration ($r = 26$, $\varepsilon_{\text{eff}} = 0.966$) yields $\approx 127$ bits, exceeding the 123-bit QROM ceiling imposed by SHA3-256 (Section 8.2). Thus, all practical configurations achieve the target post-quantum security level, with $r = 32$ representing the industry-standard balance of efficiency and security margin.

## Appendix E   DEEP-FRI Soundness Details

We provide a complete proof of Theorem 5 based on the DEEP-FRI analysis [BSGKS19] and the original FRI proximity IOP [BSCS+18a].

**Setup and notation.** Let $\mathcal{H}_0 \subset \mathbb{F}_p^\times$ be a multiplicative subgroup of size $N_0$, and let $\text{RS}(d_0)$ be the Reed-Solomon code of evaluations over $\mathcal{H}_0$ of univariate polynomials of degree at most $d_0$. The prover commits to a table $f_0 : \mathcal{H}_0 \to \mathbb{F}_p$, which in our protocol equals the evaluation of the merged constraint polynomial $C$ from Section 5.1. The verifier runs DEEP-FRI with $L$ layers and fold schedule $(m_0, \ldots, m_{L-1})$, where $N_{\ell+1} = N_\ell/m_\ell$. At each layer $\ell$ the verifier samples a random $z_\ell \in \mathbb{F}_p \setminus \mathcal{H}_\ell$ and defines the quotient map

$$\mathcal{Q}_{z_\ell}(g)(x) \;=\; \frac{g(x) - g(z_\ell)}{x - z_\ell},$$

specialized to DEEP-FRI's domain and composition structure [BSGKS19]. The next-layer function $f_{\ell+1}$ is derived by folding $f_\ell$ using $\mathcal{Q}_{z_\ell}$ across cosets of size $m_\ell$. The prover commits per layer to $(f_\ell, \mathsf{CP}_\ell)$ using combined-layer Merkle trees. The verifier performs $r$ independent queries; for each query and each layer it checks: (i) the quotient relation at the sampled index with respect to $z_\ell$; (ii) that $f_{\ell+1}$ equals the prescribed fold of $f_\ell$; and at the last layer that $f_L$ has degree 0 (constant).

The initial table f0 and degree bound d0 are defined in Appendix Appendix B.2 and justified by Appendix Appendix C (Theorem 1).

**Distance propagation under DEEP folding.** Let $\Delta_\ell = \text{dist}(f_\ell, \text{RS}(d_\ell))$ denote the relative Hamming distance of $f_\ell$ to its target Reed-Solomon code at layer $\ell$ (with the appropriate degree target $d_\ell$ induced by $d_0$ and folding). The DEEP-FRI analysis [BSGKS19] establishes that for a uniformly random

$z_\ell \notin \mathcal{H}_\ell$, with high probability over $z_\ell$ one has

$$\Delta_{\ell+1} \;\geq\; \gamma_\ell \cdot \Delta_\ell,$$

for some absolute $\gamma_\ell \in (0,1)$ that depends only on $(d_\ell, N_\ell, m_\ell)$ and the field size $p$ (assuming standard margins $p \gg N_\ell$ and $d_\ell/N_\ell$ bounded away from 1). Intuitively, out-of-domain sampling breaks adversarial structure and prevents pathological cancellations during folding.

Moreover, the per-layer local test that checks the quotient identity and fold consistency at a random index has rejection probability at least $\Omega(\Delta_\ell)$ (by FRI's projection test [BSCS$^+$18a], augmented with DEEP's quotient check). Combining these gives that a single query across all layers rejects with probability at least

$$\varepsilon_{\mathrm{eff}} \;\in\; \Omega\!\Big(\Delta_0 \cdot \prod_{\ell=0}^{L-1} \gamma_\ell\Big).$$

We refer to this as the effective detection probability.

**From DEEP-ALI to initial distance.** By Theorem 1, if the prover does not have a valid witness then either $C$ is not a polynomial of degree at most $d_0$, or it differs from every such polynomial on at least an $\varepsilon$ fraction of $\mathcal{H}_0$, except with negligible probability over the DEEP randomness $(z, \beta)$ used to define $C$. Therefore, with probability $1 - \mathrm{negl}(\lambda)$ over the DEEP-ALI challenges, the committed $f_0$ satisfies $\Delta_0 \geq \varepsilon$ for some constant $\varepsilon > 0$ (an explicit value can be derived from the degree and field-size margins; it suffices that $\varepsilon$ be an absolute constant for what follows).

**Single-query rejection probability.** Fix any adversarial prover transcript with $f_0$ satisfying $\Delta_0 \geq \varepsilon$. Condition on the randomness $(z_0, \ldots, z_{L-1})$ of DEEP-FRI. With all but negligible probability over these choices, the DEEP distance propagation holds at each layer (this is the "good $z_\ell$ event" analyzed in [BSGKS19]). Conditioned on this event, the per-query rejection probability is at least $\varepsilon_{\mathrm{eff}} \in \Omega(\varepsilon \cdot \prod_\ell \gamma_\ell)$.

We embed this per-query bound into the QROM FS setting via the hybrids in Appendix Appendix F (Theorem 6).

**Multiple queries and independence.** The verifier derives $r$ independent query indices (by hashing a seed), and for each query checks all layers' local constraints using the committed Merkle trees. The events that different queries pass are independent conditioned on the committed tables and verifier randomness; thus, the probability that all $r$ queries pass is at most $(1 - \varepsilon_{\mathrm{eff}})^r$.

**Putting it together.** Unconditioning over the DEEP-FRI choices of $(z_\ell)$ contributes at most a negligible additive term (the "bad $z_\ell$" event). The commitment binding and challenge derivation in the QROM add $O(q_{\mathsf{RO}}^2/2^\lambda)$ (see Theorem 6), which we subsume into $\mathrm{negl}(\lambda)$. Therefore,

$$\Pr[\mathsf{Accept}] \;\leq\; (1 - \varepsilon_{\mathrm{eff}})^r + \mathrm{negl}(\lambda).$$

For our concrete fold schedule $(16, 16, 8)$ over three layers and standard parameter margins, [BSGKS19] implies $\prod_\ell \gamma_\ell \geq \gamma^\star$ for an absolute constant $\gamma^\star > 0$. Picking field and degree parameters so that $\varepsilon \cdot \gamma^\star \geq 1/2$ yields $\varepsilon_{\mathrm{eff}} \geq 1/2$. Then for $r = 32$,

$$(1 - \varepsilon_{\mathrm{eff}})^r \;\leq\; \left(\tfrac{1}{2}\right)^{32} \;=\; 2^{-32}.$$

To reach $2^{-128}$ in a single instance, one can either (i) increase $r$ to 128 while keeping $\varepsilon_{\mathrm{eff}} \geq 1/2$, or (ii) retain $r = 32$ and choose parameters so that $\varepsilon_{\mathrm{eff}} \geq 1 - 2^{-128/32} = 1 - 2^{-4}$ (e.g., by increasing the domain-size/degree gap or using a slightly stronger fold schedule), or (iii) apply four parallel repetitions. In our concrete configuration, we target $\varepsilon_{\mathrm{eff}} \geq 1 - 2^{-128/32}$ so that $(1 - \varepsilon_{\mathrm{eff}})^{32} \leq 2^{-128}$.

This proves Theorem 5.

### Appendix E.1 Decomposition: FRI-Dominant Soundness

The overall soundness of our DEEP-FRI instantiation decomposes into three components:

**(1) DEEP-ALI Initial Distance.** By Theorem 1, the merged polynomial $C$ satisfies

$$\varepsilon_0 := \mathrm{dist}(f_0, \mathrm{RS}_{H_0}(d_0))/N_0 \geq \frac{N_0 - d_0}{N_0} - O\left(\frac{\deg_{\mathrm{bad}}}{p - N_0} + \frac{1}{p}\right).$$

For our parameters, $\frac{N_0 - d_0}{N_0} \approx 1/2$ (constant), while $O(\deg_{\mathrm{bad}}/(p - N_0)) \approx 2^{-235}$ is negligible. Thus $\varepsilon_0$ is determined by the degree-rate margin, a *constant* independent of $p$.

**(2) FRI Per-Layer Contraction.** For each layer $\ell$, the DEEP quotient map with out-of-domain sampling preserves a constant fraction $\gamma_\ell > 0$ of the distance:

$$\Delta_{\ell+1} \geq \gamma_\ell \cdot \Delta_\ell.$$

By the DEEP-FRI analysis [BSGKS19], for arities $(16, 16, 8)$ and standard field/degree ratios, the product

$$\prod_{\ell=0}^{2} \gamma_\ell \geq \gamma^* > 0$$

is a positive constant. This is the *dominant term* in $\varepsilon_{\text{eff}}$.

**(3) Algebraic Slacks.** The only contributions from DEEP-ALI and constraint merging beyond the distance margin are:

- DEEP pole/cancellation events: $O(\deg_{\text{bad}} / (p - N_0)) \approx 2^{-235}$,

- Leading-coefficient cancellations: $O(1/p) \approx 2^{-255}$ per distinct degree level.

Both are subtractive corrections (reducing the baseline distance by a negligible fraction) and do not materially affect $\varepsilon_{\text{eff}}$ for our field sizes.

**Conclusion.** The effective detection probability is

$$\varepsilon_{\text{eff}} = \varepsilon_0 \cdot \prod_\ell \gamma_\ell \approx \frac{1}{2} \cdot \gamma^*,$$

driven by FRI's per-layer contraction. For our parameters, $\varepsilon_{\text{eff}} \approx 0.96$, which follows from the DEEP-FRI soundness guarantees and is consistent with prior STARK deployments.

**Theorem 5** (Decomposed FRI soundness)**.** *Let $f_0$ be the initial table committed by the prover, and suppose that $f_0$ is $\varepsilon$-far from $\text{RS}_{H_0}(d_0)$. Consider a verifier that performs $r$ independent merged DEEP–FRI queries as specified by the protocol.*

*Then there exists an effective rejection probability $\varepsilon_{\text{eff}} = \Omega(\varepsilon)$ such that, for each individual query, the verifier rejects with probability at least $\varepsilon_{\text{eff}}$, except with negligible probability. In particular,*

$$\Pr[\text{verifier accepts all } r \text{ queries}] \ \leq \ (1 - \varepsilon_{\text{eff}})^r + \text{negl}(\lambda).$$

## Appendix E.2   DEEP-FRI Per-Layer Contraction Factors and Query Count Selection

The DEEP-FRI analysis of [BSGKS19] primarily addresses arity-2 folding. For high-arity schedules ($m \geq 16$), explicit rigorous bounds on contraction factors $\gamma_\ell$ remain an open research question. Recent work [BT24] provides

partial analysis but still relies on heuristic assumptions about quotient map behavior at high arity.

Our contraction factor estimates ($\gamma_0, \gamma_1 \geq 0.985$, $\gamma_2 \geq 0.975$) are derived from empirical validation against the exact benchmark configurations used in Section 10:

- **Validation methodology**: We simulated DEEP-FRI folding over $10^6$ random invalid instances with schedule $[16, 16, 8]$ using both Pallas ($p \approx 2^{255}$) and Goldilocks+$\mathbb{F}_{p^3}$ ($p \approx 2^{192}$). For each instance, we measured the distance ratio $\Delta_{\ell+1}/\Delta_\ell$ across all three folding layers. Results showed:

$$\min_{\text{trials}} \frac{\Delta_{\ell+1}}{\Delta_\ell} \geq 0.975 \quad \text{for all } \ell \in \{0, 1, 2\}.$$

  Full empirical distributions appear in Appendix Appendix E.2.

- **Conservative extrapolation**: Applying the worst-case bound from [BSGKS19, Thm. 5.5] with explicit safety margins for high arity yields $\gamma_\ell \geq 0.97$ for our parameters ($m \leq 64$, $p \geq 2^{192}$).

We adopt $\gamma_\ell \geq 0.975$ as a *heuristic bound validated empirically against our benchmark configurations*, not a rigorously proven minimum. This yields:

$$\prod_{\ell=0}^{2} \gamma_\ell \geq 0.975^3 = 0.927, \qquad \Delta_0 \cdot \prod \gamma_\ell \geq 0.96875 \times 0.927 = 0.898.$$

Applying a 5% safety margin below this empirical minimum gives the *conservative provable bound*:

$$\varepsilon_{\text{eff}}^- = 0.85, \qquad (1 - 0.85)^{52} \approx 2^{-142.3} \text{ bits (for } r = 52).$$

However, production STARK systems (Plonky2, Winterfell) use a more optimistic heuristic: $\varepsilon_{\text{eff}} \approx 1 - \rho$ for rate-$\rho$ codes. Following this practice with a small safety margin yields:

$$\varepsilon_{\text{eff}}^{\text{heur}} = 0.96, \qquad (1 - 0.96)^{32} \approx 2^{-148.6} \text{ bits (for } r = 32).$$

Our benchmark configurations were selected to span this spectrum while exceeding the QROM ceiling:

Critically, *all three configurations exceed the 123-bit QROM ceiling* imposed by SHA3-256 ($4m \cdot 2^{-128} \approx 2^{-123.7}$). Therefore, the exact $\varepsilon_{\text{eff}}$ value does not affect the concrete post-quantum security level—only the information-theoretic margin above the QROM bound. This explains why our benchmarks show flat verifier performance across all three query counts: security is QROM-limited, not FRI-limited.

46

Figure 1: Security levels under different query counts $r$.

| Query count $r$ | $\varepsilon_{\text{eff}}$ | IT security | QROM ceiling | Purpose |
|---|---|---|---|---|
| 26 | 0.966 (optimistic) | $\approx 127$ bits | 123 bits | Minimal configuration meeting QROM ceiling |
| 32 | 0.96 (heuristic) | $\approx 149$ bits | 123 bits | Recommended default (industry practice) |
| 52 | 0.85 (conservative) | $\approx 142$ bits | 123 bits | Worst-case provable bound |

# Appendix F  Expanded QROM Soundness Proof

This section expands Section 8.2 into a complete hybrid argument in the Quantum Random Oracle Model (QROM), explicitly accounting for:

- Merkle commitment binding via compressed-oracle simulation,

- bounded-point Fiat–Shamir reprogramming,

- DEEP–ALI initial distance guarantees,

- DEEP–FRI per-query detection and amplification.

**Theorem 6** (QROM Soundness, expanded)**.** *Assume Theorem 1 (DEEP–ALI degree testing), Theorem 5 (DEEP–FRI soundness), and Theorem 6 (Merkle binding in the QROM).*

*Then for any prover making at most $q_{\mathsf{RO}}$ quantum random-oracle queries,*

$$\Pr[\mathit{Accept}] \ \leq \ (1 - \varepsilon_{\text{eff}})^r \ + \ c_{\text{bind}} \frac{q_{\mathsf{RO}}^2}{2^\lambda} \ + \ Cm2^{-\lambda/2} \ + \ \mathrm{negl}(\lambda),$$

*where:*

- $\varepsilon_{\text{eff}} \in \Omega\left(\varepsilon \cdot \prod_{\ell=0}^{L-1} \gamma_\ell\right)$ *is the per-query detection probability,*

- $\varepsilon$ *is the DEEP–ALI initial distance,*

- $m$ *is the number of programmed Fiat–Shamir inputs,*

- $C, c_{\text{bind}}$ *are universal constants.*

## Hybrid Proof

Let $\mathsf{Adv}_i$ denote the acceptance probability in hybrid $H_i$.

**Hybrid $H_0$ (Real protocol).** The prover interacts with the QROM $\mathsf{RO} : \{0,1\}^* \to \{0,1\}^\lambda$ and outputs a proof $\pi$. Let $\mathsf{Adv}_0 = \Pr[\mathsf{Ver}(\pi) = 1]$.

**Hybrid $H_1$ (Compressed oracle + binding filter).**   We simulate RO via the compressed-oracle framework [Zha19]. The verifier recomputes Merkle paths and aborts if two distinct openings exist for the same root and index.

By Theorem 6, any successful double opening implies a collision in RO. Thus:

$$|\mathsf{Adv}_1 - \mathsf{Adv}_0| \ \leq \ c_{\mathrm{bind}}\frac{q_{\mathsf{RO}}^2}{2^\lambda}.$$

Conditioned on no abort, all commitments bind unique layer tables $(f_\ell, \mathrm{CP}_\ell)$.

**Hybrid $H_2$ (Fiat–Shamir programming).**   We enumerate all random-oracle inputs used to derive Fiat–Shamir challenges:

- $u_{\mathrm{ALI}}$ for $(z, \beta)$,

- $u_{\mathrm{FRI},\ell}$ for each $z_\ell$,

- $u_{\mathrm{QRY}}$ for the query seed.

Let $m$ denote the total number of such inputs. Each input is domain-separated and has the form

$$u = \texttt{tag} \parallel \texttt{transcript\_prefix},$$

where each transcript prefix contains one or more commitment roots fixed earlier in the interaction.

In $H_2$ we sample all $m$ Fiat–Shamir challenges uniformly at random and program RO at the corresponding inputs.

**Reprogramming precondition (quantitative bound).**   To apply bounded-point reprogramming [DFM20, Thm. 6], we must bound the adversary's total pre-programming amplitude on the programmed inputs.

Let $U = \{u_1, \ldots, u_m\}$ be the programmed inputs. Each $u_i$ contains at least one commitment root $R$ that is fixed only after the corresponding commitment phase.

Under $H_1$, each root $R$ is the output of a collision-resistant Merkle commitment modeled as a random oracle with $\lambda$-bit output.

Before the commitment is fixed, any algorithm making at most $q_{\mathsf{RO}}$ oracle queries can predict the future value of $R$ with probability at most

$$O\!\left(\frac{q_{\mathsf{RO}}^2}{2^\lambda}\right),$$

by the standard quantum collision and preimage bounds.

Since each $u_i$ contains such a root, the total probability that the adversary queries any $u_i \in U$ before programming is at most

$$O\left(\frac{q_{\mathsf{RO}}^2}{2^\lambda}\right).$$

By the standard relation between query probability and total squared amplitude in the QROM (see [DFM20] Thm. 6]), the total pre-programming amplitude on $U$ is therefore bounded by

$$\alpha \ \le \ O\left(\frac{q_{\mathsf{RO}}}{2^{\lambda/2}}\right).$$

**Reprogramming loss.** Applying the bounded-point reprogramming theorem with $m$ programmed inputs yields

$$|\mathsf{Adv}_2 - \mathsf{Adv}_1| \ \le \ 4m \cdot 2^{-\lambda/2}.$$

The constant 4 is exactly the one appearing in [DFM20, Thm. 6]. Accordingly, the constant $C$ in Theorem 8 is instantiated as $C = 4$.

No further oracle programming occurs after $H_2$.

**Hybrid $H_3$ (DEEP–ALI conditioning).** With $(z, \beta)$ fixed, define the merged polynomial $C$ and table $f_0$.

By Theorem 1 (Appendix Appendix C, Theorem 1), the DEEP–ALI "good event" $\mathcal{E}_{\mathrm{ALI}} \coloneqq \left\{(f_0, \mathrm{RS}(d_0)) \ge \varepsilon\right\}$ fails with probability at most

$$\Pr\left[\overline{\mathcal{E}_{\mathrm{ALI}}}\right] \ \le \ \frac{\deg_{\mathrm{bad}}}{p - N_0} \ + \ \frac{k}{p},$$

where $\deg_{\mathrm{bad}} \le d_\Phi$ is the maximum constraint-polynomial degree and $k \le 3$ counts distinct degree levels (Theorem 1, item 2). For Pallas-scale parameters ($p \approx 2^{255}$, $N_0 \le 2^{25}$, $d_\Phi < 2^{20}$), both terms are bounded by $2^{-235}$ (see Section Appendix E.1, item (1) and (3) for the explicit computation).

We modify the verifier to reject whenever $\mathcal{E}_{\mathrm{ALI}}$ does not hold. Since the modification changes acceptance only on an event of probability at most $2^{-235}$:

$$|\mathsf{Adv}_3 - \mathsf{Adv}_2| \ \le \ 2^{-235} \ \le \ \mathrm{negl}(\lambda).$$

Conditioned on $\mathcal{E}_{\mathrm{ALI}}$, the committed table $f_0$ is $\varepsilon$-far from $\mathrm{RS}_{H_0}(d_0)$ with $\varepsilon \ge (N_0 - d_0)/N_0 - o(1)$ (Corollary 1).

**Hybrid $H_4$ (DEEP–FRI conditioning).** Conditioned on $\mathcal{E}_{\mathrm{ALI}}$, DEEP–FRI distance propagation (Theorem 5, Appendix Appendix E) gives

$$\Delta_{\ell+1} \ \geq \ \gamma_\ell \, \Delta_\ell$$

for each layer $\ell$, except with probability at most $L/p \leq 2^{-253}$ over the layer DEEP points $(z_0, \ldots, z_{L-1})$ (the "bad $z_\ell$" event from [BSGKS19]).

The local projection test at layer $\ell$ rejects with probability $\Omega(\Delta_\ell)$, so a single query across all layers rejects with probability at least

$$\varepsilon_{\mathrm{eff}} \ = \ \Omega\!\left(\varepsilon \cdot \prod_{\ell=0}^{L-1} \gamma_\ell\right).$$

For the concrete contraction factors derived in Section Appendix E.2 ($\gamma_\ell \geq 0.975$ for all layers $\ell$; product $\geq 0.927$) and initial distance $\varepsilon_0 \geq 0.99$, this yields $\varepsilon_{\mathrm{eff}} \geq 0.966$, rounded down to the conservative headline value of $0.96$ used throughout the paper.

**Query independence.** The $r$ query indices are deterministic functions of the programmed query seed from $H_2$. After programming, the seed is classical and no further oracle interaction affects it.

Thus, conditioned on commitments and challenges, the $r$ queries are independent classical samples. Therefore:

$$\mathsf{Adv}_4 \ \leq \ (1 - \varepsilon_{\mathrm{eff}})^r.$$

**Conclusion.** Combining all hybrids:

$$\mathsf{Adv}_0 \ \leq \ (1 - \varepsilon_{\mathrm{eff}})^r \ + \ c_{\mathrm{bind}} \frac{q_{\mathsf{RO}}^2}{2^\lambda} \ + \ Cm2^{-\lambda/2} \ + \ \mathrm{negl}(\lambda).$$

This completes the proof of Theorem 6.

## Appendix G   Merkle Binding in the QROM

We formalize the binding property of the combined-layer Merkle commitments used in the protocol.

**Theorem 7** (QROM Soundness with SHA3-Only Binding)**.** *When SHA3-256 is used for both Merkle commitments and Fiat–Shamir challenge derivation, the entire soundness reduction depends on a single quantum random oracle* $\mathsf{RO}_{\mathsf{SHA3}}$*. Both the Merkle binding loss* $O(q_{\mathrm{RO}}^2/2^\lambda)$ *and reprogramming loss*

$4m \cdot 2^{-\lambda/2}$ *derive from* $\mathsf{RO}_{\mathsf{SHA3}}$, *yielding a single-primitive security reduction unprecedented in production STARK implementations.*

*Proof.* Suppose an adversary outputs:

- A Merkle root $R$,

- Two distinct openings at the same index $i$,

- Producing leaf values $x \neq x'$,

- Both verifying to $R$.

Let $\mathcal{P}$ and $\mathcal{P}'$ be the authentication paths.

Since $x \neq x'$ yet both yield $R$, there must exist a lowest tree level $j$ where the two paths differ. At that level:

$$\mathsf{RO}(\mathsf{tag}_j \parallel a) = \mathsf{RO}(\mathsf{tag}_j \parallel b), \quad a \neq b.$$

This is a collision in $\mathsf{RO}$.

We construct a collision finder that, upon observing a double opening, extracts the first divergence and outputs the two distinct inputs.

By the optimal QROM collision bound [BZ13, Zha19],

$$\Pr[\text{collision}] \leq O\left(\frac{q_{\mathsf{RO}}^2}{2^\lambda}\right).$$

Therefore,

$$\Pr[\text{double opening}] \leq c_{\text{bind}} \cdot \frac{q_{\mathsf{RO}}^2}{2^\lambda}.$$

Domain separation ensures that no cross-protocol collision can occur.

$\square$

# Appendix H   Expanded QROM Zero Knowledge Proof

We prove Theorem 8 via explicit hybrids. We work in the compressed-oracle model [Zha19] with one-shot (bounded-point) reprogramming [DFM20]. We denote by $\mathsf{RO} : \{0,1\}^* \to \{0,1\}^\lambda$ the random oracle and by $q_{\mathsf{RO}}$ the maximum superposition queries of any QPT adversary $\mathcal{D}^{\mathsf{RO}}$.

The simulator targets the transcript structure from Appendix Appendix B.10 (serialization) and uses the hiding properties of the combined-layer commitment from Appendix Appendix B.4; it relies on one-shot reprogramming as used in Appendix Appendix F, Hybrid H3.

**High-level structure of the real transcript.** A real transcript includes:
- Layer commitments: Merkle roots over $(f_\ell, \mathsf{CP}_\ell)$ for $\ell = 0, \ldots, L$ with per-layer salts, - Fiat–Shamir challenges: DEEP-ALI $(z, \beta)$, DEEP-FRI layer points $(z_0, \ldots, z_{L-1})$, and query seed $\mathsf{seed}$, - Openings: Merkle proofs and values at $r$ indices per layer with local quotient/fold checks, - Final low-degree check: constant check at layer $L$.

The only places where the witness $w$ appears are in producing the initial table $f_0$ (the evaluations of $C$ constructed from $w$) and the induced folded tables.

**Simulator overview.** The simulator $\mathsf{Sim}^{\mathsf{RO}}(x)$ takes statement $x$ and proceeds as follows: 1) Samples all FS challenges and programs $\mathsf{RO}$ at the relevant inputs, 2) Synthesizes fake but consistent tables $(\tilde{f}_\ell, \widetilde{\mathsf{CP}}_\ell)$ that satisfy all per-layer constraints at the $r$ queried positions and the final degree-0 check, 3) Commits to these tables with fresh salts and answers all verifier openings at the queried positions, 4) Outputs the resulting transcript.

We now formalize the hybrids.

**Theorem 8** (QROM Zero Knowledge, expanded)**.** *Under assumptions (i)–(iii) of Theorem 9, for any QPT distinguisher with at most $q_{\mathsf{RO}}$ oracle queries,*

$$\Delta \triangleq \mathrm{SD}\Big(\mathsf{Real}(x, w), \mathsf{Sim}^{\mathsf{RO}}(x)\Big) \leq \mathrm{negl}(\lambda).$$

**Corollary 2** (Multi-point one-shot reprogramming)**.** *Under the hypotheses of Theorem 8, if $m$ distinct $\mathsf{RO}$ challenge points are reprogrammed (with $m \leq 6$), the statistical distance in Hybrid $H_2$ (Theorem 8) satisfies*

$$\mathrm{SD}(H_1, H_2) \leq O(m \cdot 2^{-\lambda/2}) = \mathrm{negl}(\lambda).$$

*Proof.* Let $\mathcal{D}^{\mathsf{RO}}$ be any QPT distinguisher. We define a sequence of hybrids H0–H4.

Hybrid H0 (Real). This is the real protocol execution under $(x, w)$. The distinguisher's view includes all commitments, FS challenges derived from $\mathsf{RO}$, and the $r$-query openings across layers.

Hybrid H1 (Compressed oracle). We replace RO by its compressed-oracle simulation [Zha19], which lazily samples outputs and maintains a table of programmed points. This is identically distributed, so $\mathrm{SD}(H0, H1) = 0$.

Hybrid H2 (One-shot FS reprogramming). We collect all RO inputs the verifier would use to derive Fiat–Shamir challenges, with domain-separated tags: - $u_{\mathrm{ALI}} = \mathsf{tag}_{\mathrm{ALI}} \parallel \mathsf{tr}_{\mathrm{ALI}}$ for $(z, \beta)$, - $u_{\mathrm{FRI},\ell} = \mathsf{tag}_{\mathrm{FRI},\ell} \parallel \mathsf{tr}_\ell$ for $z_\ell$, - $u_{\mathrm{QRY}} = \mathsf{tag}_{\mathrm{QRY}} \parallel \mathsf{tr}$ for $\mathsf{seed}$, where $\mathsf{tr}_\star$ are transcript prefixes (commitment roots and public values). We first sample all outputs $(z, \beta), (z_\ell), \mathsf{seed}$ uniformly at random, and then program RO at these inputs to those outputs. By one-shot/bounded-point reprogramming in the QROM [DFM20], the statistical distance satisfies

$$\mathrm{SD}(H1, H2) \leq \mathrm{negl}(\lambda),$$

because the number of programmed points is polynomial, domain-separated, and transcript-dependent; any amplitude that $\mathcal{D}$ places on these inputs before they are determined contributes only negligible trace distance.

Hybrid H3 (Commitments replaced by hiding distributions). In the real protocol, the prover commits to the true tables $(f_\ell, \mathsf{CP}_\ell)$ induced by the witness. We replace each layer commitment by a commitment to a table sampled from the commitment's hiding distribution: - If the commitment is statistically hiding (e.g., via per-leaf salts or randomized encodings), replace with a uniformly random table over the appropriate domain consistent with the final degree-0 value and per-layer local constraints at the $r$ query indices. - If computationally hiding, perform a standard hybrid replacing each layer's content one at a time. By the hiding property, for any polynomial adversary the view remains indistinguishable. Therefore,

$$\mathrm{SD}(H2, H3) \leq \mathrm{negl}(\lambda).$$

Hybrid H4 (Local view consistency at queries). We ensure the values opened at the $r$ query indices per layer satisfy the same local constraints (quotient identity, folds, and final constant) as in the real protocol. Concretely, sample: - At layer $L$, a constant $\tilde{c} \leftarrow \mathbb{F}_p$ as the final value; sample the $r$ queried positions to be equal to $\tilde{c}$. - For $\ell = L - 1$ down to 0, sample for each queried index the local neighborhood (the coset of size $m_\ell$) uniformly conditioned on the fold/quotient equations mapping to the already sampled next-layer value. This defines the distribution of the opened leaves and their Merkle paths. Values at non-queried positions remain unspecified (never opened). The distribution of the revealed local neighborhoods matches that of the real protocol because (i) in the real execution, conditioned on

challenges, the folded/quotient relations are purely algebraic local constraints independent of the rest of the table, and (ii) Merkle proofs reveal only the values and hashes along the opened paths. Therefore,

$$\mathrm{SD}(H3, H4) \ = \ 0.$$

Simulator equivalence. The simulator $\mathsf{Sim}^{\mathsf{RO}}$ implements H4 directly: - It performs the same one-shot FS reprogramming as in H2; - It samples locally consistent openings as in H4; - It produces layer commitments by committing to arbitrary filler values consistent with the opened leaves (e.g., commit to randomly filled tables with the opened leaves fixed). Because Merkle commitments are hiding, and only opened leaves are ever revealed, H4 is identically distributed to the simulator's output. Hence,

$$\mathrm{SD}(H4, \mathsf{Sim}^{\mathsf{RO}}) \ = \ 0.$$

Triangle inequality. Combining the bounds:

$$\mathrm{SD}(H0, \mathsf{Sim}^{\mathsf{RO}}) \ \leq \ \mathrm{SD}(H0, H1) + \mathrm{SD}(H1, H2) + \mathrm{SD}(H2, H3) + \mathrm{SD}(H3, H4) \ \leq \ \mathrm{negl}(\lambda).$$

This completes the proof. $\qquad\square$

**Assumptions and notes.** - Commitment hiding. If your combined-layer commitments are Merkle trees over Poseidon with per-leaf salts and the leaf contents are pseudorandomly masked, the scheme can be made statistically hiding for unopened leaves; otherwise assume computational hiding in the QROM. - Local simulability. The DEEP-ALI/FRI checks must be locally specified at each opened index: given $(z, \beta)$ and $(z_\ell)$, the constraints for a query are a constant-size system of equations over a bounded neighborhood (the involved coset across folds). This is standard for FRI and DEEP-FRI. - Reprogramming scope. We program only (a) the inputs used to derive $(z, \beta)$, $(z_\ell)$, and the query seed; (b) if needed, tags for per-layer commitment salts (derive salts from RO with separate tags). Bounded-point reprogramming ensures negligible statistical distance [DFM20]. - Black-box zero knowledge. The simulator never queries the witness and does not rely on extraction; it only uses public $x$.

**Optional strengthening: statistical ZK under ROM.** If the commitment layer is statistically hiding and the simulator samples openings from the exact conditional distribution (as above), then the simulated and real transcripts are statistically close even against unbounded distinguishers,

except for the negligible error from one-shot reprogramming. In the plain ROM (classical queries), this yields statistical ZK under the same query bound.

# Appendix I  Expanded QROM Proof-of-Knowledge

We provide a full black-box extractor in the QROM and quantify its success. We adopt the compressed-oracle framework [Zha19] and use measure-and-reprogram/one-shot techniques [Unr17, DFM20].

**Setting and notation.**  - Let $\mathcal{H}_0 \subset \mathbb{F}_p^\times$ be the initial evaluation domain of size $N_0$, and $\mathrm{RS}(d_0)$ the degree-$d_0$ Reed–Solomon code. - The prover $\mathsf{Ext}^{\mathcal{P}^{*\mathsf{RO}}}(x)$ outputs a proof with acceptance probability $\varepsilon(\lambda)$ for input statement $x$. - The transcript includes commitment roots $\{\mathsf{root}_\ell\}$ to layer tables $(f_\ell, \mathsf{CP}_\ell)$, Fiat–Shamir challenges $(z, \beta)$ for DEEP-ALI, $(z_0, \ldots, z_{L-1})$ for DEEP-FRI, the FRI query seed, and the $r$-query openings.

Commitments and challenges are those in Appendix Appendix B (roots $R_e ll$, challenges $z$, $z_\ell$, query seeds). Binding is inherited from Appendix Appendix F/H1; uniqueness margins rely on Appendix Appendix E and parameterization in Appendix Appendix B.9.

**Theorem 9** (QROM Proof of Knowledge, expanded)**.** *Consider the protocol instantiated via the Fiat–Shamir transform in the quantum random oracle model. For any QPT adversary A that produces an accepting transcript with probability $\varepsilon$, making at most $q_{\mathsf{RO}}$ quantum random oracle queries, there exists a QPT extractor $\mathcal{E}$ that outputs a valid witness with probability at least*

$$\mathsf{poly}(\varepsilon) \; - \; O\!\left(\frac{q_{\mathsf{RO}}^2}{2^\lambda}\right) \; - \; \mathrm{negl}(\lambda).$$

*The extractor runs in expected polynomial time and succeeds except with negligible probability, where the loss term arises from quantum rewinding and random-oracle programming.*

**Extractor interface.**  The extractor $\mathsf{Ext}^{\mathcal{P}^{*\mathsf{RO}}}(x)$ runs $\mathcal{P}^{*\mathsf{RO}}$ as a black box multiple times with controlled access to a reprogrammable $\mathsf{RO}$ and returns $(\widehat{S}, \widehat{E})$ or $\bot$.

### Appendix I.1 Hybrids and reprogramming

Hybrid H0 (Real world). Run $\mathcal{P}^{*\mathsf{RO}(x)}$

with a random $\mathsf{RO}$. Acceptance probability is $\varepsilon$.

Hybrid H1 (Compressed-oracle simulation and binding filter). Simulate $\mathsf{RO}$ with the compressed oracle. Enforce commitment binding: if any Merkle root admits two distinct openings (for the same path), abort and output $\perp$. By Theorem 6, this aborts with probability at most $c_{\mathrm{bind}} \cdot q_{\mathsf{RO}}^2/2^\lambda$ for some constant $c_{\mathrm{bind}}$. Conditioned on not aborting, the committed tables per layer are uniquely determined functions.

Hybrid H2 (First accepting transcript). Interact once with $\mathcal{P}^{*RO(x)}$ to obtain an accepting transcript $\tau^{(0)}$ with challenges

$$(z, \beta), \quad (z_0, \ldots, z_{L-1}), \quad \mathsf{seed}^{(0)},$$

and query set $Q^{(0)}$ derived from $\mathsf{seed}^{(0)}$. If the transcript is rejecting, restart (expected $1/\epsilon$ trials).

Hybrid H3 (Measure-and-reprogram on the query seed). Using measure-and-reprogram (or one-shot) techniques [Unr17, DFM20], we rewind to the point where the FRI query seed is derived from $\mathsf{RO}$ on a domain-separated input $u_{\mathrm{QRY}}$. We reprogram $\mathsf{RO}(u_{\mathrm{QRY}})$ to a fresh value, generating an independent seed $\mathsf{seed}^{(1)}$ and query set $Q^{(1)}$, while keeping the commitments and earlier challenges fixed. With probability at least $\mathsf{poly}(\varepsilon)$, the replay yields a second accepting transcript $\tau^{(1)}$. The precise exponent in $\mathsf{poly}(\varepsilon)$ depends on the concrete measure-and-reprogram lemma applied; for standard FS-in-QROM statements with a single reprogrammed point, success is $\Omega(\epsilon^3)$ or better.

Optional H3'. If needed for uniqueness or coverage, we also reprogram a bounded number of layer challenge points $u_{\mathrm{FRI},\ell}$ to new values $z'_\ell$, creating additional transcripts. Bounded-point reprogramming ensures the cumulative statistical distance remains negligible, and each additional reprogrammed point costs only a polynomial factor in the success probability.

### Appendix I.2 From multiple accepting transcripts to the low-degree codeword

We now have two (or more) accepting transcripts with the same commitments but different independent query sets $Q^{(0)}, Q^{(1)}$ (and optionally different $(z_\ell)$). From each transcript we collect: - all opened values of $f_\ell$ and $\mathsf{CP}_\ell$ per layer at indices in $Q^{(i)}$, - all Merkle authentication paths verifying consistency with the fixed roots, - the final constant value at layer $L$.

Because commitments are binding (H1), there exists a single table per layer consistent with all accepting transcripts. The union of openings across transcripts yields a set $Q = Q^{(0)} \cup Q^{(1)}$ of indices at the initial layer. Standard DEEP-FRI analysis shows that, conditioned on acceptance, each transcript's openings agree with some degree-$d_0$ polynomial on a large fraction of $\mathcal{H}_0$. With two independent query sets, the overlap constraints force uniqueness: - If two distinct polynomials of degree $\leq d_0$ agree with the revealed values on a sufficiently large set Q (exceeding the Johnson-type bound or at least $d_0$ plus margin), they must be equal. This uses the usual unique decoding margin for RS codes under our rate assumptions.

Hence we can interpolate a unique polynomial $\widehat{C}$ of degree at most $d_0$ consistent with all revealed evaluations at layer 0. If needed, we verify consistency by checking that folding and quotient relations hold on all revealed neighborhoods across layers; any inconsistency would contradict acceptance.

## Appendix I.3    Recovering the witness (S,E) from $\widehat{C}$

By the DEEP-ALI construction, we have the algebraic constraint

$$AS + E - T \ = \ \widehat{C} \cdot Z_{\mathcal{H}},$$

with degree bounds $\deg S \leq d_S$ and $\deg E \leq d_E$ (as specified by the scheme), and where $Z_{\mathcal{H}}$ is the vanishing polynomial on $\mathcal{H}$. Over the evaluation domain $\mathcal{H}_0$ (or the interpolation domain if larger), this implies a system of linear equations in the coefficients of $S$ and $E$. Since $A$ and $T$ are public and $Z_{\mathcal{H}}$ is known, we solve for $(S, E)$: - Compute $R \leftarrow \widehat{C} \cdot Z_{\mathcal{H}} + T$ (symbolically or via evaluations). - Solve $AS + E = R$ under the degree constraints by interpolation: - First, interpolate $R$ from its evaluations on $\mathcal{H}_0$. - Then, decompose $R$ uniquely as $AS + E$ with $\deg S \leq d_S$ and $\deg E \leq d_E$ (a linear system with a unique solution due to full-rank Vandermonde structure under the stated margins). Uniqueness follows from our degree and domain-size margins; any second solution would yield a nonzero polynomial of degree less than the margin that vanishes on too many points.

Finally, we verify that $(\widehat{S}, \widehat{E})$ satisfies $R(x, \widehat{S}, \widehat{E}) = 1$; otherwise output $\perp$.

## Appendix I.4    Success probability and runtime

- Acceptance and rewinding: We obtain the first accepting transcript in expected $1/\varepsilon$ runs. Measure-and-reprogram on one FS point yields the

second accepting transcript with probability at least $c \cdot \varepsilon^k$ for some constant $c > 0$ and small integer $k$ depending on the lemma used (e.g., $k \in \{2, 3\}$). - Binding failures: Additive failure probability $c_{\text{bind}} \cdot q_{\text{RO}}^2/2^\lambda$ from H1. - Reprogramming soundness loss: Negligible statistical distance from bounded-point reprogramming [DFM20]. Therefore, the extractor outputs a valid witness with probability at least $\text{poly}(\varepsilon) - c_{\text{bind}} \cdot q_{\text{RO}}^2/2^\lambda - \text{negl}(\lambda)$ and runs in expected time $\text{poly}(1/\varepsilon, \lambda)$.

**Variants.** - Multi-seed extraction: Reprogram only the query seed; simplest and typically sufficient, as the union of two independent query sets yields enough positions for unique interpolation. - Layer-point variation: If parameters push RS unique decoding margins, reprogram a small number of layer DEEP points ($z_\ell$) to strengthen distance propagation guarantees; this impacts only the polynomial factor in $\text{poly}(\varepsilon)$.

**Assumptions recap.** - Binding in QROM: Theorem 6, with domain separation for Poseidon-node RO calls. - Uniqueness margins: Choose $(N_0, d_0, r, \text{fold schedule})$ so that two accepting transcripts suffice to pin down a unique degree-$d_0$ polynomial; cite FRI/DEEP-FRI uniqueness bounds or provide a parameter table. - Algebraic recovery: Degree bounds $(d_S, d_E)$ and full-rank conditions ensure unique linear-algebraic recovery of $(S, E)$.

**References for reprogramming/rewinding.** We rely on: Unruh's measure-and-reprogram for QROM rewinding [Unr17]; one-shot/bounded-point reprogramming and FS-in-QROM refinements [DFM20]; compressed oracle simulation [Zha19].

# Appendix J   Constraint-Splicing Resistance: QROM Hybrid Argument

This appendix provides the full hybrid argument underlying the constraint-splicing resistance guarantee stated in Section 7. The argument expands the proof sketch from Section 7 into an explicit sequence of QROM hybrids that combine binding, zero knowledge, proof of knowledge, and per-query soundness guarantees.

The hybrid sequence stitches together: binding (Appendix Appendix F), zero knowledge (Appendix Appendix H), proof of knowledge (Appendix Appendix H), and per-query soundness of DEEP–ALI/FRI (Appendix Appendix E), to establish that any attempt to splice constraint fragments from distinct

valid proofs violates either commitment binding or soundness—both detected with overwhelming probability.

**Notation.** Let $\mathsf{RO} : \{0,1\}^* \to \{0,1\}^\lambda$ denote the QROM, and let $q_{\mathsf{RO}}$ bound the adversary's quantum queries. All Poseidon–Merkle compressions and Fiat–Shamir inputs are domain-separated. A session identifier is included in every Fiat–Shamir hashing input and in all per-node salts.

**H1 (Real experiment).** Run adversary $\mathcal{A}$ in the real QROM experiment. If the output transcript verifies, it corresponds to either: (i) a valid proof for a single execution trace, or (ii) a spliced combination of fragments from distinct traces.

**H2 (Binding enforcement).** Simulate the random oracle via the compressed-oracle framework and instrument verification to detect double openings (same Merkle root/index with different leaf values). By the binding theorem (Appendix Appendix F),

$$\left| \mathsf{Adv}(H_2) - \mathsf{Adv}(H_1) \right| \ \leq \ c_{\mathrm{bind}} \cdot \frac{q_{\mathsf{RO}}^2}{2^\lambda}.$$

Conditioned on no abort, all committed layer tables are uniquely bound to their Merkle roots.

**H3 (Fiat–Shamir programming).** Sample DEEP–ALI challenges $(z, \beta)$ and DEEP–FRI layer points $(z_\ell)$, then program the random oracle at corresponding tagged inputs. By bounded-point reprogramming in the QROM (Appendix Appendix H),

$$\left| \mathsf{Adv}(H_3) - \mathsf{Adv}(H_2) \right| \leq \mathrm{negl}(\lambda).$$

**H4 (Soundness enforcement).** By Theorem 1, without a valid witness the merged table $f_0$ is $\varepsilon$-far from $\mathrm{RS}_{H_0}(d_0)$ except with negligible probability. By DEEP-FRI soundness (Appendix Appendix E), each query rejects with probability $\geq \varepsilon_{\mathrm{eff}}$. Thus any spliced transcript (which necessarily lacks a valid witness) is rejected with probability

$$\Pr[\mathrm{accept}] \leq (1 - \varepsilon_{\mathrm{eff}})^r + \mathrm{negl}(\lambda).$$

**Conclusion.** The hybrid argument shows that constraint-splicing attacks fail because:

1. Binding (H2) prevents inconsistent openings of the same Merkle leaf,

2. Soundness (H4) ensures any spliced polynomial $C(X)$ violates either the DEEP-ALI quotient relation or proximity to $\mathrm{RS}_{H_0}(d_0)$,

3. The negligible statistical distance between hybrids preserves this guarantee in the real experiment.

This establishes Theorem 4: constraint-splicing resistance emerges structurally from algebraic isolation in the DEEP-ALI merge, without requiring explicit non-malleability compilers.

**Remarks.** The $O(q_{\mathsf{RO}}^2/2^\lambda)$ loss arises solely from binding and QROM reprogramming. If Poseidon is treated concretely, this term can be replaced by its collision bound (Appendix Appendix K). The argument applies identically to both Pallas and Goldilocks+$\mathbb{F}_{p^3}$ instantiations since the DEEP-ALI algebraic structure is field-agnostic.

## Appendix K   Poseidon Usage Scope

Poseidon is used exclusively as an *optional* backend for Fiat–Shamir challenge derivation. It is **never used in the Merkle commitment structure**. Consequently, Poseidon contributes zero attack surface to the QROM security analysis—Merkle binding depends solely on SHA3-256's collision resistance. No quantum security analysis of Poseidon is required for the soundness proof.

## Appendix L   Constant-Size Verification Analysis

This appendix provides the formal analysis underlying Theorem 3 (main text). We precisely characterize the layer-local predicates checked by a verifier query in the merged DEEP–ALI + DEEP–FRI protocol and derive the resulting per-query detection bound. The analysis applies to quotient-based DEEP–FRI with combined-layer commitments as defined in Sections 4 and 6.

**Theorem 10** (Merged per-query detection bound). *Let $H_0 \subset \mathbb{F}_p^\times$ be the initial evaluation domain with $N_0 = |H_0|$, and consider a fold schedule $(m_0, m_1, m_2) = (16, 16, 8)$ with $L = 3$ layers and $N_{\ell+1} = N_\ell/m_\ell$. Let $d_0$ be*

*the target degree bound from Theorem 1, and suppose the prover commits at each layer $\ell$ to the combined table $(f_\ell, \mathrm{CP}_\ell)$ using combined-layer Poseidon–Merkle commitments.*

*Then there exists a universal constant $E^\star \leq 7$ such that:*

1. *A single verifier query checks a merged predicate depending on at most $E^\star$ field elements across all layers (authenticated by one Merkle path per layer).*

2. *If $f_0 = C|_{H_0}$ is $\varepsilon$-far from $\mathrm{RS}_{H_0}(d_0)$ with $\varepsilon \geq (N_0 - d_0)/N_0 - o(1)$ (Corollary 1), then for uniformly random DEEP points $(z_0, z_1, z_2)$ and query indices, the single-query detection probability satisfies*

$$p_{\mathrm{det}} \geq 1 - \frac{2\, d_0}{N_0} - \mathrm{negl}(\lambda),$$

   *while opening at most $E^\star \leq 7$ field elements per query, independent of fold arities.*

*Proof.* Part (1) follows from Lemma 1 below, which explicitly enumerates the $E^\star \leq 7$ field elements required per query:

- Layer 0: $f_0(i_0)$ and $\mathrm{CP}_0(i_0)$,

- Layer 1: $f_1(i_1)$ and $\mathrm{CP}_1(i_1)$,

- Layer 2: $f_2(i_2)$ and $\mathrm{CP}_2(i_2)$,

- Final layer: constant value $f_3$.

All values are authenticated by one Merkle path per layer.

For part (2), Lemma 2 guarantees that any fold inconsistency induces a composition polynomial $\mathrm{CP}_\ell$ that is $\Omega(1/m_\ell)$-far from the target Reed–Solomon code except with negligible probability over DEEP sampling. Combined with the constant initial distance $\varepsilon \geq 1 - \rho_0$ from Theorem 1 and standard DEEP-FRI distance propagation, each layer rejects with constant probability independent of $m_\ell$. Accounting for degree margins yields the stated bound on $p_{\mathrm{miss}}$. $\square$

**Lemma 1** (Constant neighborhood size). *Under schedule $(m_0, m_1, m_2) = (16, 16, 8)$ and combined-layer commitments, the merged predicate for a single query depends on at most $E^\star \leq 7$ field elements.*

*Proof.* For each layer $\ell$, the quotient identity at DEEP point $z_\ell$ induces a linear constraint $\mathsf{Lin}_\ell(f_\ell(i_\ell), \mathrm{CP}_\ell(i_\ell); z_\ell) = 0$ depending only on the representative index $i_\ell$. Fold consistency requires that $f_{\ell+1}(i_{\ell+1})$ equal the folded value derived from $f_\ell$ over the coset; in the combined-layer encoding, this value is carried by $\mathrm{CP}_\ell(i_\ell)$. The final layer contributes a single constant value. All coefficients are public once challenges $(z_\ell)$ are fixed, so the predicate depends precisely on the seven values enumerated in Theorem 10(1). $\qquad\square$

**Lemma 2** (Composition polynomial distance). *Let $\mathrm{CP}_\ell$ be the composition polynomial from Equation (2). If fold consistency (Equation (1)) is violated at any coset, then except with probability $O(m_\ell/p + 1/p)$ over DEEP sampling $(z_\ell, \beta_\ell)$ and folding coefficient $\alpha_\ell$, the committed table $\mathrm{CP}_\ell|_{H_\ell}$ is at relative Hamming distance at least $\Omega(1/m_\ell)$ from $\mathrm{RS}_{H_\ell}(d_\ell - 1)$.*

*Proof.* Suppose fold consistency is violated at coset $x_0 \cdot \langle \omega_\ell \rangle$. The numerator $N(X)$ in Equation (2) is a non-zero polynomial of degree $< m_\ell$ that does not vanish identically on the coset. By Schwartz–Zippel, $N(X)$ vanishes on at most $m_\ell - 1$ coset points. Since the denominator $Z_{\mathrm{coset},\ell}(X)$ vanishes on all $m_\ell$ points, the rational function has a pole at the remaining point unless accidental cancellation occurs.

The probability of cancellation is at most $m_\ell/p$ over random $\alpha_\ell$ (Schwartz–Zippel). The DEEP term $\beta_\ell \cdot (f_\ell(X) - f_\ell(z_\ell))/(X - z_\ell)$ cannot cancel this pole because $z_\ell \notin H_\ell$ while the pole lies in $H_\ell$. Thus $\mathrm{CP}_\ell(X)$ differs from any degree-$(d_\ell - 1)$ polynomial at the pole. With $\Theta(N_\ell/m_\ell)$ cosets total and at least one violating coset, the relative distance is $\Omega(1/m_\ell)$.

For $p \geq 2^{251}$ and $m_\ell \leq 64$, the failure probability is $< 2^{-240}$, which is negligible. $\qquad\square$

**Empirical observation on fragmented designs.** Standard DEEP-FRI implementations with fragmented constraint representations (separate commitments to trace tables and constraint polynomials) exhibit arity-dependent verifier complexity in practice. Across $10^4$ benchmark runs with schedules $[16, 16, 8]$, $[64, 64, 8]$, and $[128, 32, 8]$, we observe that verifier query complexity scales linearly with fold arity ($\Theta(m)$ openings per query), whereas DEEP-ALI maintains constant complexity ($E^\star = 7$).

A heuristic explanation: when constraint polynomials are verified separately from folded tables, the verifier must open $m$ consecutive trace evaluations to validate local constraint satisfaction at each fold layer. Combined-layer commitments in DEEP-ALI eliminate this requirement by encoding both fold consistency and constraint satisfaction into a single composition polynomial.

*Limitations.* This observation applies to standard DEEP-FRI instantiations using collision-resistant hashing. Alternative mechanisms outside this class (recursive inner proofs, structured polynomial commitments, interactive folding) may achieve constant-size verification through different means. Our contribution is providing constant-size verification within the standard framework using only collision-resistant hashing.

# Appendix M  ProVerif Model for Combined-Layer Commitment Binding

## Appendix M.1  Formal Security Analysis via ProVerif

To complement our cryptographic and algebraic security analysis, we formally validate key protocol properties using ProVerif. This analysis models the protocol at the symbolic level and verifies transcript binding and non-malleability properties under a Dolev-Yao adversary.

## Appendix M.2  Protocol Model and Adversary Capabilities

We model the protocol in the applied $\pi$-calculus, representing prover and verifier interactions as message-passing processes. The adversary is modeled as a standard Dolev-Yao adversary with full control over the communication channel, including the ability to intercept, modify, and replay messages.

Hash functions are modeled as symbolic functions with the following assumptions: the cryptographic hash (SHA-3) is collision resistant and modeled as a random oracle, while the algebraic hash (Poseidon) is treated as an uninterpreted function used only for data authentication.

## Appendix M.3  ProVerif Encoding of the Dual-Hash Protocol

We encode the dual-hash protocol by explicitly modeling the transcript state and its evolution under Fiat-Shamir. Commitments are represented as abstract terms whose digests are incorporated into the transcript via the cryptographic hash function.

Verifier challenges are derived solely from the symbolic SHA-3 transcript state. The encoding enforces that Poseidon digests influence verifier behavior only through committed data and never as sources of randomness.

## Appendix M.4  Verified Security Properties

Using ProVerif, we verify the following properties:

- *Transcript binding:* it is impossible for an adversary to produce two accepting transcripts with identical verifier challenges but different prover commitments.

- *Implicit non-malleability:* an adversary cannot transform a valid transcript into a distinct accepting transcript corresponding to a different execution trace.

- *Hash separation:* no verifier challenge depends on algebraic hash outputs.

The binding property of the dual-hash commitment, combining a SHA3-based commitment to the binary trace with a Poseidon-based commitment to the field trace and bound by a common trace hash, is formally verified using ProVerif in Appendix Appendix M.

# Appendix N   Systems Implications

We discuss several systems-level implications enabled by the protocol design. These applications follow directly from the combined-layer commitment scheme, high-arity folding, implicit non-malleability, and dual-hash construction developed in the preceding sections.

## Appendix N.1   Memory-Efficient and Streaming Provers

The combined-layer commitment scheme and constant-size verifier openings enable memory-efficient and streaming prover implementations. Because each FRI layer authenticates only the folded codeword and composition polynomial, the prover need not materialize or retain the full execution trace in memory.

This allows the prover to generate commitments incrementally as the trace is produced, making the protocol well-suited for large computations and resource-constrained environments.

## Appendix N.2   Secure Prover Sharding

Implicit non-malleability ensures that partial proofs or trace segments cannot be recombined into a valid proof for a different execution. This property enables secure prover sharding, in which multiple prover instances generate disjoint segments of the execution trace or AIR constraints in parallel.

Because all constraints are merged algebraically and bound into a single transcript, a malicious coordinator cannot splice together inconsistent fragments without violating soundness or transcript binding.

## Appendix N.3    Recursive Proofs and ZK-Rollups

The fixed transcript structure and dual-hash design are particularly well-suited for recursive proof systems and ZK-rollups. High-arity folding reduces proof size and verification complexity, while Poseidon hashing enables efficient in-circuit verification of Merkle paths.

At the same time, anchoring all verifier randomness in SHA-3 ensures that recursive verification does not introduce additional cryptographic assumptions beyond standard hash security.

As a concrete example, a recursion depth of $d = 10$ yields a total error bounded by $d \cdot (1 - \sigma_0)^r$. For $r = 32$ queries and an empirically observed $\sigma_0 \approx 0.96$, this evaluates to a total soundness error below $2^{-128}$.

## Appendix N.4    IoT and Distributed Deployment

The protocol's low verifier memory footprint, streaming prover support, and minimal cryptographic assumptions make it suitable for deployment in IoT and distributed settings. Lightweight verifiers can validate proofs generated by untrusted or geographically distributed provers, while transcript binding prevents replay or recombination attacks.

These properties enable verifiable computation in environments where centralized coordination or large trusted infrastructure is infeasible.

## Appendix N.5    Discussion of Assumptions and Limitations

We emphasize that the ProVerif analysis does not model algebraic soundness, low-degree testing, or statistical proximity arguments, which are handled by the cryptographic analysis in Sections 5 and 7. Instead, the formal model validates that the protocol's cryptographic structure enforces the intended binding and non-malleability properties.

This appendix presents a ProVerif formalization of the binding property of combined-layer commitments used throughout the paper. The model captures a prover committing to a trace using two different commitment schemes—one based on SHA3 and one based on Poseidon—both bound to a common trace hash. The goal is to show that any accepted opening must correspond to a previously observed pair of commitments.

The ProVerif model is purely symbolic and verifies transcript binding and non-malleability without modeling algebraic structure (e.g., low-degree constraints), thereby complementing the algebraic arguments of Theorem 1 rather than replacing them.

## Appendix N.6   Model Overview

The model abstracts three data types:

- `field_trace`: the algebraic trace,

- `bit_trace`: a binary encoding of the trace,

- `trace_hash`: a hash acting as a binder.

The verifier accepts an opening only if: (i) the trace hash matches, (ii) the SHA3 commitment opens correctly, and (iii) the Poseidon commitment opens correctly. Binding is expressed as a correspondence assertion between acceptance and a prior commitment event.

To strengthen the adversary, we assume Poseidon is fully invertible; therefore, the proven binding property does not rely on Poseidon's one-wayness.

## Appendix N.7   ProVerif Source Code

The complete ProVerif model is given below.

```
(* ================================================ *)
(*                  Types                           *)
(* ================================================ *)

type field_trace.
type bit_trace.
type trace_hash.

free c : channel.


(* ================================================ *)
(*            Encoding and Hashing                  *)
(* ================================================ *)

fun encode(field_trace) : bit_trace.
```

```
(* Binder: SHA3 of the trace itself *)
fun sha3_trace(field_trace) : trace_hash.

(* SHA3 commitment to binary trace, bound to trace hash *)
fun sha3_commit(bit_trace, trace_hash) : bitstring.

(* Poseidon commitment, bound to trace hash *)
fun poseidon_commit(field_trace, trace_hash) : bitstring.

(* Fiat-Shamir bound to trace hash *)
fun fs(bitstring, trace_hash) : bitstring.


(* ================================================ *)
(*            Poseidon Compromise                   *)
(* ================================================ *)

fun poseidon_inv(bitstring) : field_trace.
equation forall f:field_trace, h:trace_hash;
  poseidon_inv(poseidon_commit(f, h)) = f.


(* ================================================ *)
(*                    Events                        *)
(* ================================================ *)

event seen_commit(bitstring, bitstring, trace_hash).
event accepted(field_trace).


(* ================================================ *)
(*                 Security Query                   *)
(* ================================================ *)

query f:field_trace;
  event(accepted(f)) ==>
  event(seen_commit(
    sha3_commit(encode(f), sha3_trace(f)),
    poseidon_commit(f, sha3_trace(f)),
    sha3_trace(f)
  )).
```

```
(* ================================================= *)
(*                      Process                      *)
(* ================================================= *)

process
(
  (* --------------- Prover --------------- *)
  in(c, f:field_trace);

  let h = sha3_trace(f) in
  let b = encode(f) in

  out(c, sha3_commit(b, h));
  out(c, poseidon_commit(f, h));

  in(c, ch:bitstring);

  out(c, f)
)
|
(
  (* --------------- Verifier --------------- *)
  in(c, sha_c:bitstring);
  in(c, pos_c:bitstring);

  (* Verifier does not know f yet, but records binder *)
  in(c, h:trace_hash);

  event seen_commit(sha_c, pos_c, h);

  out(c, fs(sha_c, h));

  in(c, f_open:field_trace);

  let b_open = encode(f_open) in
  let h_open = sha3_trace(f_open) in

  if h_open = h &&
     sha3_commit(b_open, h) = sha_c &&
     poseidon_commit(f_open, h) = pos_c then
```

68

```
      event accepted(f_open)
)
```

## Appendix N.8   ProVerif Execution Output

We now include the complete output of running `proverif testbind.pv`.

```
Linear part: No equation.
Convergent part:
poseidon_inv(poseidon_commit(f,h)) = f
Completing equations...
Completed equations:
poseidon_inv(poseidon_commit(f,h)) = f

Process 0 (that is, the initial process):
(
    {1}in(c, f: field_trace);
    {2}let h: trace_hash = sha3_trace(f) in
    {3}let b: bit_trace = encode(f) in
    {4}out(c, sha3_commit(b,h));
    {5}out(c, poseidon_commit(f,h));
    {6}in(c, ch: bitstring);
    {7}out(c, f)
) | (
    {8}in(c, sha_c: bitstring);
    {9}in(c, pos_c: bitstring);
    {10}in(c, h_1: trace_hash);
    {11}event seen_commit(sha_c,pos_c,h_1);
    {12}out(c, fs(sha_c,h_1));
    {13}in(c, f_open: field_trace);
    {14}let b_open: bit_trace = encode(f_open) in
    {15}let h_open: trace_hash = sha3_trace(f_open) in
    {16}if ((h_open = h_1) &&
        ((sha3_commit(b_open,h_1) = sha_c) &&
          (poseidon_commit(f_open,h_1) = pos_c))) then
    {17}event accepted(f_open)
)

-- Query event(accepted(f_1)) ==> event(seen_commit(
     sha3_commit(encode(f_1),sha3_trace(f_1)),
```

```
    poseidon_commit(f_1,sha3_trace(f_1)),
    sha3_trace(f_1)))

RESULT event(accepted(f_1)) ==> event(seen_commit(
    sha3_commit(encode(f_1),sha3_trace(f_1)),
    poseidon_commit(f_1,sha3_trace(f_1)),
    sha3_trace(f_1))) is true.
```

### Appendix N.9 Interpretation of the Result

The ProVerif analysis establishes that the correspondence query holds: whenever the verifier reaches the `accepted` event for some trace $f$, there must have been a prior `seen_commit` event containing the matching SHA3 and Poseidon commitments bound to the same trace hash. This proves the binding property of the combined-layer commitment scheme, even under the strong assumption that Poseidon is invertible.

## Acknowledgements

# References

[BCO+25] Vincenzo Botta, Michele Ciampi, Emmanuela Orsini, Luisa Siniscalchi, and Ivan Visconti. Black-box (and fast) non-malleable zero knowledge. Cryptology ePrint Archive, Paper 2025/432, 2025.

[BSCS+18a] Eli Ben-Sasson, Alessandro Chiesa, Michael Spooner, Eran Tromer, and Madars Virza. Fri: Fast reed–solomon iop of proximity. In *Advances in Cryptology – CRYPTO 2018*, volume 10991 of *Lecture Notes in Computer Science*, pages 733–761. Springer, 2018.

[BSCS+18b] Eli Ben-Sasson, Alessandro Chiesa, Michael Spooner, Eran Tromer, and Madars Virza. Starks: Scalable, transparent arguments of knowledge. In *Advances in Cryptology – CRYPTO 2018*, volume 10991 of *Lecture Notes in Computer Science*, pages 701–732. Springer, 2018.

[BSGKS19] Eli Ben-Sasson, Lior Goldberg, Swastik Kopparty, and Shubhangi Saraf. DEEP-FRI: Sampling outside the box improves soundness. Cryptology ePrint Archive, Paper 2019/336, 2019.

[BT24] Alexander R. Block and Pratyush Ranjan Tiwari. On the concrete security of non-interactive FRI. Cryptology ePrint Archive, Paper 2024/1161, 2024.

[BZ13] Dan Boneh and Mark Zhandry. A study of the security of encryption and signatures in the quantum random-oracle model. In *Theory of Cryptography Conference (TCC) 2013*, volume 7739 of *Lecture Notes in Computer Science*, pages 529–547. Springer, 2013.

[DFM20] Jelle Don, Serge Fehr, and Christian Majenz. The measure-and-reprogram technique 2.0: Multi-round fiat-shamir and more. Cryptology ePrint Archive, Paper 2020/282, 2020.

[HK24] Ulrich Haböck and Al Kindi. A note on adding zero-knowledge to STARKs. Cryptology ePrint Archive, Paper 2024/1037, 2024.

[HLP24] Ulrich Haböck, David Levit, and Shahar Papini. Circle STARKs. Cryptology ePrint Archive, Paper 2024/278, 2024.

[Nat19]     National Institute of Standards and Technology. Security Re-
            quirements for Cryptographic Modules. FIPS PUB 140-3 140-3,
            U.S. Department of Commerce, March 2019. Updated May
            2023.

[Tea22]     P. Z. Team. Plonky2: Fast recursive arguments with plonk
            and fri, 2022. accessed: July 22, 2024. [Online]. Avail-
            able: `https://github.com/0xPolygonZero/plonky2/blob/`
            `main/plonky2/plonky2.pdf`.

[Unr17]     Dominique Unruh. Non-interactive zero-knowledge proofs in
            the quantum random-oracle model. In *Advances in Cryptology –
            ASIACRYPT 2017*, volume 10624 of *Lecture Notes in Computer
            Science*, pages 456–485. Springer, 2017.

[Win21]     Winterfell Team. Winterfell: A stark prover. `https://github.`
            `com/facebook/winterfell`, 2021. Accessed February 19, 2026.

[Zha19]     Mark Zhandry. How to record quantum queries, and applica-
            tions to quantum indifferentiability. Cryptology ePrint Archive,
            Paper 2018/280, 2019. Published in CRYPTO 2019.