# Adversarial Technical Review

*Header-Only Verification and the State Commitment Invariant (v3)*

A Deliberately Hostile Examination of Boundary Conditions

## Threat Model and Review Scope

**This review evaluates the header-only verification invariant under adversarial conditions.** In scope: cryptographic correctness of the invariant given honest majority consensus, binding commitments, and available headers. The model guarantees that accepted proofs are consistent with consensus-finalized state commitments—not that finalized commitments reflect semantically valid protocol execution. Out of scope: header availability (eclipse attacks, network partitions), prover incentives, liveness guarantees, and application-specific semantic validation. Assumptions: honest consensus majority (>50% or >2/3 depending on mechanism), collision-resistant hash functions, binding commitment schemes, deterministic state transition rules, and existence of at least one honest prover at query time. This review tests where these assumptions break and what additional guarantees would require stronger premises.

## 0. Review Objective and Methodology

The paper claims that header-only verification with state commitments and bounded proofs guarantees "correctness under standard consensus safety assumptions." The core model consists of:

- Headers $H_i$ at height i
- Canonical state $S_i$ after applying all valid transitions up to height i
- Commitment $C_i = Commit(S_i)$ embedded in headers
- Proofs $\pi_x$ for state queries x

The paper assumes binding commitments, deterministic canonical encoding, and prefix-monotonic chain selection.

This is a strong invariant. The adversarial examination focuses on forcing explicit boundaries between what the model guarantees and what it does not, identifying where the term "correctness" becomes ambiguous, and demonstrating attacks that succeed without violating the core invariant.

## 1. Critical Ambiguity: What "Correctness" Means

The invariant states:

> *"A verifier need not execute transitions if consensus commits to a binding representation of post-transition state."*

This guarantees consistency with finalized commitments. It does not guarantee semantic validity of state transitions unless consensus participants already enforced protocol rules correctly.

### Hostile Objection (Critical)

**The model proves *"consensus agreed on $C_i$,"* not *"$C_i$ reflects a valid state transition."***

If consensus finalizes a commitment to a state that violates protocol semantics—whether due to implementation bugs, governance overrides, or adversarial majority—the header-only verifier accepts it by design. The verifier has no mechanism to detect semantic invalidity because it does not execute transitions.

The paper acknowledges this in Section 13.1 ("What Header-Only Verification Does Not Do") but the abstract and invariant statements use "correctness" without qualification, creating ambiguity.

### Required Revision

Replace "guarantee correctness" with:

> *"Guarantee state query correctness and consistency relative to consensus-finalized commitments, assuming consensus participants enforce protocol semantics."*

# 2. Commitment Scheme: Canonical Encoding Underspecified

Section 2.1 requires commitments be binding, succinct, verifiable, and based on "canonical encoding."

## Hostile Objection (High)

**"Canonical encoding" is not a cryptographic primitive—it is an implementation requirement that determines whether independent nodes compute identical $C_i$ for logically equivalent states.**

Cross-client commitment divergence is a known failure mode in multi-implementation systems. If clients disagree on:

- Key ordering (lexicographic vs insertion order)
- Serialization format (big-endian vs little-endian, padding rules)
- Domain separation tags
- Trie structure (Patricia vs binary)

then honest full nodes will produce different $C_i$ for the same logical state, causing consensus failure or chain splits. This is an architectural limitation, not a violation of the invariant, but it must be explicitly acknowledged.

### Required Revision

Add to Section 2.1: "Canonical encoding must be fully specified with reference test vectors. Implementation disagreement on serialization, key ordering, or trie construction violates the determinism assumption and is out of scope for this invariant."

# 3. Chain Selection: Cryptographic vs Operational Guarantees

Section 2.2 assumes prefix-monotonic chain selection: Weight(C') > Weight(C) for extensions.

## Hostile Objection (High)

The cryptographic invariant holds, but the light client's ability to identify the canonical chain depends on correct header delivery—a network-layer problem outside the model's scope.

The paper correctly notes this in Section 9.3 ("Header Availability Model"), but earlier sections blur the distinction between:

- **Ledger integrity given correct headers (in scope)**
- **Ability to obtain correct headers (out of scope)**

**Required Revision**

Abstract and Section 4 should explicitly state: "This invariant guarantees correctness of verification given access to valid headers. Header sourcing, eclipse resistance, and anti-censorship mechanisms are orthogonal operational requirements not addressed by this model."

# 4. Adversarial Boundary Testing

The following attacks succeed without violating the core invariant. Each exploits a documented limitation rather than a cryptographic weakness. These attacks illustrate information-theoretic boundaries—not implementation bugs—that would require stronger assumptions to prevent.

## Attack 1: Consensus Finalizes Semantically Invalid State

**Scenario:** An adversarial majority or implementation bug causes consensus to finalize $C_i = Commit(S_i)$ where $S_i$ violates protocol rules (e.g., double-spend, invalid signature, incorrect balance calculation).

**Verifier behavior:** Accepts $C_i$ as valid because the header is consensus-signed and the commitment is binding.

**Outcome:** Attack succeeds.

*Interpretation: Not a violation of the invariant. The model guarantees consistency with finalized commitments, not enforcement of protocol semantics. Preventing this attack requires trusted execution or assuming honest consensus enforces semantic validity—an assumption outside the model's scope.*

*This is a boundary condition, not a contradiction.*

## Attack 2: Proof Withholding (Liveness Denial)

**Scenario:** Adversary controls all reachable provers and withholds proofs $\pi_x$ for state queries.

**Verifier behavior:** Cannot obtain witnesses for balance checks, transaction inclusion, or other queries.

**Outcome:** Attack succeeds (liveness failure).

*Interpretation: The paper acknowledges this in Section 9.2: proof withholding causes liveness failure but does not enable forgery. The adversary cannot produce false proofs that pass verification. This is an availability attack, not an integrity violation. Mitigation requires prover redundancy, caching, or incentive mechanisms—all out of scope.*

*This exploits a documented limitation: the model assumes at least one honest prover is available.*

## Attack 3: Eclipse the Header Feed

**Scenario:** Adversary controls the victim's network view and presents an alternative header chain with lower cumulative weight but higher local appearance of validity.

**Verifier behavior:** Accepts the presented chain as canonical if it satisfies local validity checks (linkage, consensus rules, commitment consistency).

**Outcome:** Attack succeeds if the victim cannot access honest headers.

*Interpretation: This is the classic SPV failure mode. The paper explicitly notes in Section 9.3 that header availability and eclipse resistance are out of scope. The invariant guarantees integrity of accepted data, not the ability to obtain correct data.*

*Preventing this requires bonded relayers, multi-source aggregation, or checkpoint authorities —operational solutions outside the cryptographic model.*

## Attack 4: Operational Complexity Degrades Practical Security

**Scenario:** Deployed system introduces proof markets, caching layers, and P2P routing. An adversary exploits operational dependencies to degrade availability via targeted censorship or resource exhaustion.

**Verifier behavior:** Experiences degraded service quality but cannot detect whether failures are malicious or benign.

**Outcome:** Attack succeeds (practical security degradation without integrity violation).

*Interpretation: The paper's claim that security is "non-amplifying" (min(S_headers, S_consensus)) applies to integrity of accepted proofs, not to operational robustness. Real deployments add attack surface through infrastructure dependencies.*

*The non-amplifying claim should be qualified: "for integrity of cryptographically verified data, assuming correct header delivery."*

# 5. Strengths to Preserve

Despite boundary vulnerabilities, the paper's core contributions are sound:

- **The invariant is cleanly stated and quotable.** Section 3 provides a precise formalization that is architecturally agnostic.
- **Non-normative stance is appropriate.** The paper correctly avoids prescribing specific cryptographic primitives or consensus mechanisms.
- **Section 13 (Non-Goals) is well-scoped.** The limitations are documented. The abstract and summary should reference this section more prominently to prevent misinterpretation.

# 6. Verdict and Required Revisions

## ACCEPT WITH MINOR REVISIONS

The cryptographic framework is internally sound. The identified limitations are information-theoretic or architectural boundaries inherent to header-only verification, not implementation defects or logical contradictions. The following revisions will eliminate reviewer ambiguity:

1. **Redefine "correctness" as consensus-relative.** Abstract and Section 3 should state: "This invariant guarantees state query correctness and commitment consistency relative to consensus-finalized state, assuming consensus participants

enforce protocol semantics. It does not verify that finalized states satisfy semantic validity conditions."

2. **Specify commitment determinism requirements.** Section 2.1 should add: "Canonical encoding must be defined by a complete specification with reference test vectors. Cross-client divergence on serialization, key ordering, or trie construction violates the determinism assumption and causes consensus failure independent of this invariant."

3. **Qualify the non-amplifying security claim.** Section 11 should state: "The min(S_headers, S_consensus) bound applies to integrity of accepted proofs given correct header delivery. Operational infrastructure (proof markets, relayers, caching) may introduce availability vulnerabilities that are orthogonal to this cryptographic guarantee."

# Summary of Adversarial Findings

This hostile review validates the core invariant under adversarial stress-testing while identifying three critical areas requiring explicit clarification:

- **Scope of correctness guarantees:** The model guarantees consistency with finalized commitments, not enforcement of semantic validity. Consensus-finalized invalid states are accepted by design.
- **Commitment determinism:** Canonical encoding requires complete specification. Implementation disagreement causes consensus failure orthogonal to verification correctness.
- **Integrity vs availability separation:** The security bound applies to accepted data assuming correct header delivery. Operational robustness is out of scope.

All identified attacks exploit documented limitations without violating the invariant. No logical contradictions exist. The model is sound within its explicitly stated assumptions.

With the specified revisions, this work provides a rigorous foundation for header-only verification systems with unambiguous boundary conditions.

*— End of Adversarial Technical Review —*

December 2025

*Prepared for publication as companion critique to research note v3*