

# Security Analysis of DAG-BFT Consensus: A Formal Comparison with Nakamoto Consensus

Q-NarwhalKnight Protocol Team

December 17, 2025

## Abstract

We present a formal security analysis of the Q-NarwhalKnight DAG-BFT consensus protocol, comparing its security guarantees with Bitcoin’s Nakamoto consensus. We prove safety and liveness properties under explicit assumptions, quantify attack costs under economic models, and identify irreducible tradeoffs between the two approaches. Our analysis demonstrates that DAG-BFT provides deterministic finality with approximately 300x energy reduction per unit of economic security, while accepting a lower Byzantine threshold (33% vs 51%) and requiring weak subjectivity for node bootstrapping. We explicitly characterize failure modes, separate cryptographic from economic and social guarantees, and acknowledge limitations that remain open problems.

## 1 Introduction

Distributed consensus protocols face fundamental tradeoffs between security, efficiency, and trust assumptions. Bitcoin’s Nakamoto consensus [1] achieves probabilistic finality through Proof-of-Work (PoW), externalizing security costs into energy expenditure. Byzantine Fault Tolerant (BFT) protocols achieve deterministic finality through voting, internalizing security costs into capital at risk.

This paper provides a formal analysis of Q-NarwhalKnight, a DAG-structured BFT protocol, comparing its security properties with Nakamoto consensus. We do not claim universal superiority; rather, we characterize the precise tradeoffs each system accepts.

### 1.1 Contributions

1. Formal safety and liveness proofs for DAG-BFT under partial synchrony
2. Quantitative attack cost comparison between PoW and PoS models
3. Explicit characterization of synchrony failure modes
4. Separation of cryptographic, economic, and social security layers
5. Honest assessment of irreducible tradeoffs and open problems

## 2 System Model

### 2.1 Network Model

**Assumption 2.1** (Partial Synchrony). *There exists a Global Stabilization Time (GST) and a message delay bound  $\Delta$  such that after GST, all messages between correct processes are delivered within  $\Delta$  time units.*

**Assumption 2.2** (Byzantine Threshold). *At most  $f < n/3$  of  $n$  validators are Byzantine, where Byzantine validators may behave arbitrarily (including colluding, equivocating, or remaining silent).*

**Assumption 2.3** (Cryptographic Primitives). *The following primitives are secure:*

- *Digital signatures (SQIsign or Dilithium5) are existentially unforgeable under chosen-message attacks*
- *Hash functions (SHA3-256, BLAKE3) are collision-resistant*
- *VDF evaluation requires sequential computation (no parallel speedup)*

## 2.2 Validator Model

**Definition 2.4** (Validator). *A validator  $v_i$  is a tuple  $(pk_i, sk_i, stake_i)$  where:*

- *$pk_i, sk_i$  are public/private key pairs for signing*
- *$stake_i \geq stake_{min}$  is bonded collateral subject to slashing*

**Definition 2.5** (Voting Power). *Validator  $v_i$  has voting power  $w_i = \min(stake_i, cap)$  where  $cap$  is a quadratic voting cap to limit plutocracy.*

## 3 DAG-BFT Protocol

### 3.1 Block Structure

**Definition 3.1** (DAG Block). *A block  $B$  is a tuple  $(h, parents, txs, \sigma, vdf\_proof)$  where:*

- *$h = H(parents || txs || vdf\_proof)$  is the block hash*
- *$parents \subseteq \{h_1, \dots, h_k\}$  references  $\geq 2$  parent blocks*
- *$txs$  is an ordered list of transactions*
- *$\sigma = Sign(sk_i, h)$  is the proposer's signature*
- *$vd\_proof$  proves sequential computation for leader election*

### 3.2 Finality Rule

**Definition 3.2** (Finality). *A block  $B$  is finalized when it is included in the causal history of blocks signed by validators with combined voting power  $\geq 2f + 1$ .*

---

**Algorithm 1** DAG-BFT Finality Check

---

```

1: function ISFINALIZED( $B, DAG$ )
2:    $supporters \leftarrow \emptyset$ 
3:   for each block  $B'$  in  $DAG$  where  $B \in ancestors(B')$  do
4:      $supporters \leftarrow supporters \cup \{proposer(B')\}$ 
5:   end for
6:    $weight \leftarrow \sum_{v \in supporters} w_v$  return  $weight \geq 2f + 1$ 
7: end function

```

---

## 4 Security Analysis

### 4.1 Safety

**Theorem 4.1** (Safety). *Under Assumptions 2.2 and 2.3, no two conflicting blocks can both be finalized.*

*Proof.* Assume for contradiction that conflicting blocks  $B$  and  $B'$  (where neither is an ancestor of the other) are both finalized.

By Definition 3.3, finalization of  $B$  requires supporters with total weight  $\geq 2f + 1$ . Similarly for  $B'$ .

Total validator weight is  $3f + 1$  (assuming  $n = 3f + 1$  for tightness).

If both  $B$  and  $B'$  are finalized:

$$|\text{supporters}(B)| \geq 2f + 1 \quad (1)$$

$$|\text{supporters}(B')| \geq 2f + 1 \quad (2)$$

By pigeonhole, the intersection:

$$|\text{supporters}(B) \cap \text{supporters}(B')| \geq (2f + 1) + (2f + 1) - (3f + 1) = f + 1 \quad (3)$$

At least  $f + 1$  validators supported both conflicting blocks. Since at most  $f$  are Byzantine, at least one honest validator signed blocks supporting both  $B$  and  $B'$ .

But honest validators only sign blocks extending their current view. If an honest validator supported  $B$ , they would only later support  $B'$  if  $B \in \text{ancestors}(B')$  or vice versa—contradicting our assumption that  $B$  and  $B'$  conflict.

Therefore, no two conflicting blocks can both be finalized.  $\square$   $\square$

### 4.2 Liveness

**Theorem 4.2** (Liveness). *Under Assumptions 2.1, 2.2, and 2.3, any transaction submitted by a correct client is eventually finalized.*

*Proof.* After GST, messages are delivered within  $\Delta$ .

A correct validator receiving a valid transaction will include it in their next proposed block.

With  $\geq 2f + 1$  correct validators, at least one correct validator will propose a block containing the transaction within bounded time.

Correct validators extend the DAG by referencing recent blocks. Within  $O(\Delta)$  time, the block containing the transaction will be referenced by  $\geq 2f + 1$  validators' blocks.

By Definition 3.3, the block is finalized.  $\square$   $\square$

**Corollary 4.3** (Finality Time). *Under partial synchrony, finality is achieved within  $O(\Delta)$  time after GST, where  $\Delta$  is the message delay bound.*

### 4.3 Synchrony Failure Modes

**Proposition 4.4** (Fail-Stop Behavior). *Under network partition or asynchrony, DAG-BFT halts (no new finality) rather than producing conflicting finalized blocks.*

*Proof.* Safety (Theorem 4.1) holds unconditionally—it does not depend on synchrony assumptions.

Liveness (Theorem 4.2) requires partial synchrony. Under partition:

- Partitioned validators cannot communicate
- Neither partition has  $\geq 2f + 1$  honest validators (assuming balanced partition)

- No new blocks achieve finality

The system exhibits fail-stop behavior: it halts rather than producing incorrect state. This is preferable to fail-unsafe behavior where conflicting state could be finalized.  $\square$   $\square$

## 5 Attack Cost Analysis

### 5.1 Nakamoto Consensus (PoW)

**Definition 5.1** (51% Attack Cost (PoW)). *The cost to execute a double-spend attack in Nakamoto consensus is:*

$$C_{PoW} = \int_0^T (\text{hash\_rate}_{\text{attack}} \cdot \text{energy\_cost}) dt + \text{hardware\_cost} \quad (4)$$

where  $T$  is the attack duration and  $\text{hash\_rate}_{\text{attack}} > 0.5 \cdot \text{hash\_rate}_{\text{total}}$ .

**Remark 5.2.** *PoW attack cost is operational—the attacker retains hardware after the attack. Hash rate can also be rented (e.g., NiceHash), reducing capital requirements.*

### 5.2 DAG-BFT (PoS)

**Definition 5.3** (33% Attack Cost (PoS)). *The cost to violate safety in DAG-BFT is:*

$$C_{PoS} = \int_0^S P(s) \cdot \frac{dS}{ds} ds + \text{slashing\_loss} \quad (5)$$

where  $S = 0.33 \cdot \text{total\_stake}$ ,  $P(s)$  is the price at stake level  $s$ , and  $\text{slashing\_loss} = S$  (100% slashing for equivocation).

**Lemma 5.4** (Market Impact). *Acquiring stake fraction  $\alpha$  of circulating supply incurs superlinear cost due to price impact:*

$$\text{Cost}(\alpha) = \int_0^\alpha P_0 \cdot e^{\lambda s} ds = \frac{P_0}{\lambda} (e^{\lambda \alpha} - 1) \quad (6)$$

where  $\lambda$  is the market depth parameter and  $P_0$  is initial price.

**Theorem 5.5** (Attack Cost Comparison). *Under typical market conditions, the net attack cost for DAG-BFT is comparable to or higher than PoW despite the lower threshold, due to:*

1. Market impact during stake acquisition
2. Irreversible stake destruction via slashing
3. Inability to “return” stake (unlike hash rate rental)

*Proof.* Let Bitcoin market cap be  $M_{BTC}$  and Q-NarwhalKnight market cap be  $M_{QNK}$ .

For Bitcoin 51% attack via rental:

$$C_{BTC} \approx 0.51 \cdot \text{hash\_rate} \cdot \text{rental\_rate} \cdot \text{duration} \quad (7)$$

Estimated: \$500M–\$1B for sustained attack.

For Q-NarwhalKnight 33% attack:

$$C_{QNK} = \text{Cost}(0.33) + 0.33 \cdot M_{QNK} \cdot \text{staking\_ratio} \quad (8)$$

With  $\lambda \approx 3$  (moderate liquidity), acquiring 33% incurs  $\approx 10\times$  price impact.

If  $M_{QNK} = \$1B$  and staking ratio = 60%:

$$Naive\_cost = 0.33 \cdot 0.6 \cdot \$1B = \$198M \quad (9)$$

$$With\_slippage \approx \$500M - \$2B \quad (10)$$

$$Plus\_slashing = total\_loss \quad (11)$$

The attacker loses their entire stake upon detection (slashing), making the attack cost equal to acquisition cost—which is *not recoverable*.  $\square$   $\square$

### 5.3 Secondary Market Effects

**Remark 5.6** (Liquidity Risks). *The attack cost model assumes direct market acquisition. Secondary markets may reduce costs:*

- **Liquid staking derivatives:** Tradeable staked tokens reduce market impact
- **Lending markets:** Borrowed stake enables attacks without capital commitment
- **Validator cartels:** Coordinated validators share risk

*These are active areas of research and represent evolving threat vectors.*

## 6 Comparative Analysis

### 6.1 Energy Efficiency

**Theorem 6.1** (Energy Comparison). *DAG-BFT requires approximately  $300\times$  less energy per unit of economic security compared to Nakamoto consensus.*

*Proof.* Define energy per unit security as:

$$E_{sec} = \frac{Annual\_Energy}{Market\_Cap\_Security} \quad (12)$$

For Bitcoin:

$$E_{BTC} = \frac{150 \text{ TWh/year}}{\$1T} = 150 \text{ GWh}/\$B \quad (13)$$

For Q-NarwhalKnight (assuming 1000 full nodes at 100W):

$$E_{QNK} = \frac{1000 \cdot 100W \cdot 8760h}{\$1B} \quad (14)$$

$$= \frac{0.876 \text{ TWh/year}}{\$1B} \approx 0.5 \text{ GWh}/\$B \quad (15)$$

Ratio:  $\frac{150}{0.5} = 300\times$ .  $\square$   $\square$

**Remark 6.2.** *This comparison assumes stake-based security is economically equivalent to hash-based security—a contested assumption (see Section 7).*

### 6.2 Finality

**Proposition 6.3** (Finality Comparison). *Nakamoto consensus provides probabilistic finality with security exponentially increasing in confirmation depth. DAG-BFT provides deterministic finality within bounded time.*

Property	Nakamoto (PoW)	DAG-BFT (PoS)
Finality type	Probabilistic	Deterministic
Time to finality	~60 min (6 blocks)	<3 seconds
Reorg possible	Yes (with sufficient hash)	No (after finality)

Table 1: Finality comparison

### 6.3 Weak Subjectivity

**Definition 6.4** (Weak Subjectivity). *A protocol exhibits weak subjectivity if new nodes cannot securely determine the canonical chain from genesis without trusting an external checkpoint.*

**Proposition 6.5** (Weak Subjectivity in PoS). *DAG-BFT requires weak subjectivity; Nakamoto consensus does not.*

*Proof.* In Nakamoto consensus, the canonical chain is the one with most cumulative work. A new node can verify this from genesis by checking PoW proofs—no external trust required.

In DAG-BFT, the canonical chain depends on stake distribution, which is not self-verifiable from blockchain data alone. An attacker with historical keys could create an alternate history. New nodes must trust a recent checkpoint to distinguish real from fake chains.  $\square$   $\square$

**Remark 6.6.** *Weak subjectivity is a fundamental trust assumption in PoS, not a bug to be fixed. It represents a tradeoff: faster finality and lower energy in exchange for checkpoint trust.*

## 7 Security Guarantee Layers

We explicitly separate three guarantee types:

### 7.1 Layer 1: Cryptographic (Strongest)

These guarantees hold unconditionally given cryptographic assumptions:

- Signature unforgeability (Assumption 2.3)
- Hash collision resistance
- VDF sequentiality

No social coordination required. Failure requires breaking cryptographic primitives.

### 7.2 Layer 2: Economic (Conditional)

These guarantees assume rational actors and functioning markets:

- Attack cost exceeds potential gain
- Slashing deters equivocation
- Staking yield maintains validator participation

Can fail under: irrational attackers, market manipulation, cartelization.

### 7.3 Layer 3: Social (Weakest)

These guarantees rely on human coordination:

- Checkpoint distribution and trust
- Governance response to attacks
- Client implementation diversity

**Remark 7.1.** *Bitcoin also has social-layer dependencies (UASF, social consensus on chain validity). Neither system is purely trustless.*

## 8 Open Problems and Limitations

### 8.1 Acknowledged Limitations

1. **Lower Byzantine threshold:** 33% vs 51% is a real tradeoff
2. **Weak subjectivity:** Checkpoint trust is fundamental, not eliminable
3. **DAG complexity:** Higher implementation and verification burden
4. **Governance risks:** On-chain governance has known failure modes
5. **Single client:** Currently Rust-only implementation

### 8.2 Open Research Questions

1. Can recursive SNARKs reduce weak subjectivity burden?
2. What are optimal slashing parameters under various market conditions?
3. How do liquid staking derivatives affect attack economics?
4. What formal verification approaches scale to DAG protocols?

## 9 Conclusion

We have presented a formal security analysis of DAG-BFT consensus, proving safety and liveness under explicit assumptions, and comparing attack costs with Nakamoto consensus.

**Key findings:**

1. DAG-BFT provides deterministic finality with  $\sim 300\times$  energy reduction
2. Attack costs are comparable despite lower threshold (due to slashing and market impact)
3. Weak subjectivity is a fundamental tradeoff, not a flaw
4. Both systems have social-layer dependencies

**The choice between systems is domain-specific:**

- Nakamoto consensus optimizes for self-verifiable history and minimal trust assumptions
- DAG-BFT optimizes for fast finality, energy efficiency, and smart contract composability

Neither is universally superior. The appropriate choice depends on which tradeoffs are acceptable for the target application.

## References

- [1] Nakamoto, S. (2008). Bitcoin: A Peer-to-Peer Electronic Cash System.
- [2] Castro, M., & Liskov, B. (1999). Practical Byzantine Fault Tolerance.
- [3] Buchman, E., Kwon, J., & Milosevic, Z. (2018). The latest gossip on BFT consensus.
- [4] Danezis, G., et al. (2022). Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus.
- [5] Buterin, V., et al. (2020). Combining GHOST and Casper.