

# Consensus Mechanisms

## Lesson 4: How Blockchains Agree on Truth

Prof. Joerg Osterrieder

Spring 2026

# Learning Objectives

After this lesson, you will be able to:

- Explain the Byzantine Generals Problem
- Describe how Proof of Work achieves consensus
- Compare Proof of Work and Proof of Stake
- Evaluate trade-offs between different consensus mechanisms

**Prerequisites:** Lessons 1-3 (Intro, Blockchain, Cryptography)

---

is how decentralized networks agree without central authority

Cons

- 1 The Consensus Problem
- 2 Proof of Work
- 3 Proof of Stake
- 4 Other Mechanisms
- 5 Choosing a Consensus Mechanism

## The Consensus Problem

## The Problem:

- Distributed network with no central authority
- Nodes may be offline, slow, or malicious
- Need to agree on single version of truth

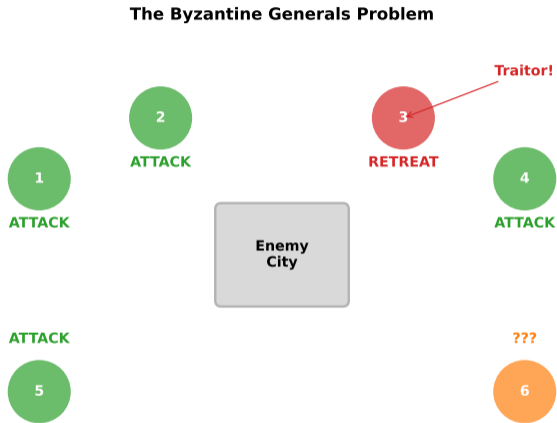
## Requirements for Consensus:

- **Agreement:** All honest nodes decide same value
- **Validity:** Decided value was proposed by some node
- **Termination:** All honest nodes eventually decide
- **Fault tolerance:** Works even if some nodes fail

---

consensus, blockchains cannot function

With



*Problem: How can loyal generals agree when some may be traitors?*

can generals coordinate when some may be traitors sending conflicting messages? Named after the Byzantine Empire, known for court intrigue and untrustworthy generals.

How

## The Problem (1982):

- Generals must agree to attack or retreat
- Messages can be lost or corrupted
- Some generals may be traitors

## The Solution:

- Need at least  $3f + 1$  generals to tolerate  $f$  traitors
- *Why  $3f+1$ ?* With  $f$  traitors, you need  $2f + 1$  honest nodes for majority. Since traitors may send conflicting messages, you need extra  $f$  to ensure honest majority prevails.
- Multiple rounds of message passing
- Majority voting with verification

**Blockchain Innovation:** Bitcoin solved this for open networks using economics—making identity creation costly through Proof of Work. Creating fake identities (Sybil attack) requires real computational resources, replacing the need for a known participant list.

---

BFT requires knowing all participants; Bitcoin makes creating fake identities prohibitively expensive

Class

## Proof of Work

# What is Proof of Work?

**Definition:** Consensus mechanism where participants prove computational work to add blocks.

## The Mining Puzzle:

- Find nonce such that:  $H(\text{block header}) < \text{target}$
- *(A hash is a large number in hexadecimal; “less than target” means the hash’s numeric value is below the threshold.)*
- Hash must start with certain number of zeros
- Only way to find: brute-force trial and error

## Why It Works:

- Hard to find valid nonce (requires work)
- Easy to verify (just compute one hash)
- Cannot cheat without doing the work

---

converts electricity into security

PoW

## Step-by-Step:

- 1 Collect pending transactions from mempool
- 2 Create block header with Merkle root
- 3 Try nonce = 0, compute hash
- 4 If hash  $>$  target, increment nonce and retry
- 5 If hash  $<$  target, broadcast winning block
- 6 Other miners verify and accept
- 7 Winner receives block reward + fees

## Bitcoin Stats:

- Target: Adjust every 2016 blocks ( 2 weeks)
- Goal: One block every 10 minutes

---

miners try billions of nonces per second

Bitco

# Difficulty Adjustment

**Problem:** Network hash rate changes over time.

**Solution:** Automatic difficulty adjustment

- If blocks too fast: increase difficulty
- If blocks too slow: decrease difficulty
- Bitcoin: Adjusts every 2016 blocks

**Formula:**

$$\text{New Difficulty} = \text{Old Difficulty} \times \frac{\text{Actual Time}}{\text{Expected Time}}$$

---

has increased by  $10^{14}$  since Bitcoin's launch

Diffic

## 51% Attack:

- Attacker needs majority of hash power
- Can rewrite recent blocks
- Cannot steal funds or create coins
- Can double-spend their own transactions

## Economic Defense:

- Cost: Billions in hardware and electricity
- Reward: Potentially double-spend some transactions
- Risk: Devalue the network you're attacking

---

Bitcoin is economically irrational

Attac

**Sybil Attack:** When an adversary creates many fake identities to gain disproportionate influence. Named after a case study of dissociative identity disorder.

### Strengths:

- Proven security (15+ years for Bitcoin)
- Permissionless participation
- Objective chain selection (most work)
- Sybil resistant

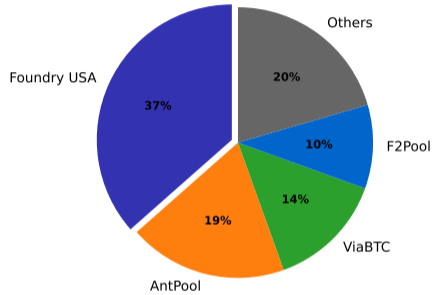
### Weaknesses:

- High energy consumption
- Hardware centralization (ASICs—Application-Specific Integrated Circuits)
- Slow finality
- Limited throughput

---

prioritizes security and decentralization over efficiency

**Bitcoin Mining Pool Distribution (as of January 2026)**



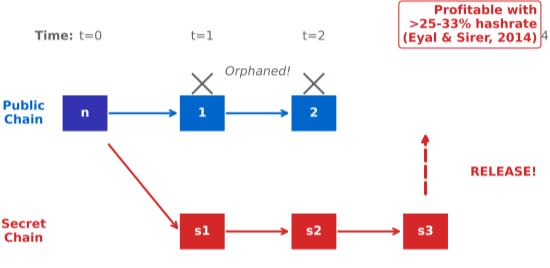
**Top 4 pools control ~80% of hashrate**

---

concentration creates censorship and 51% attack risks

Pool

## Selfish Mining Attack



**Attack Strategy:**

1. Mine secretly, don't broadcast
2. Wait until ahead by 1+ blocks
3. Release longer chain, orphaning honest blocks

mining can be profitable with 25-33% of hash power—less than 51%

## Strategy:

- Mine blocks but don't broadcast immediately
- Build a secret longer chain
- Release when ahead to orphan honest blocks

## Why It Works:

- Honest miners waste work on orphaned blocks
- Attacker gets more than proportional rewards
- Effective with  $\gamma$  (network connectivity advantage;  $\gamma$  represents the attacker's network connectivity advantage)

## Defense Mechanisms:

- Uniform tie-breaking (Bitcoin doesn't use this)
- Shorter block times make attack harder
- Monitoring for suspicious mining patterns

---

described by Eyal & Sirer (2013)—remains a theoretical concern

# Proof of Stake

# What is Proof of Stake?

**Definition:** Consensus mechanism where validators stake tokens as collateral.

**Key Idea:**

- Replace computational work with economic stake
- Validators selected based on staked amount
- Malicious behavior = lose staked tokens (slashing)

**Why It Works:**

- Validators have “skin in the game”
- Attacking hurts your own investment
- Economic incentives align with network security

---

converts capital at risk into security

PoS

## Methods:

- **Random selection:** Weighted by stake
- **Coin age:** Stake  $\times$  time held
- **Delegation:** Users delegate to validators

## Ethereum (Post-Merge):

- Minimum stake: 32 ETH
- Random selection per slot (12 seconds)
- Committees of validators vote
- Rewards for correct attestations

---

prevents prediction and manipulation

Rand

# Slashing: The Punishment Mechanism

## What Gets Slashed:

- Double signing (proposing two blocks)
- Surround voting (conflicting attestations)
- Being offline (minor penalty)

## Consequences:

- Lose portion of staked tokens
- Forced exit from validator set
- Cannot rejoin immediately

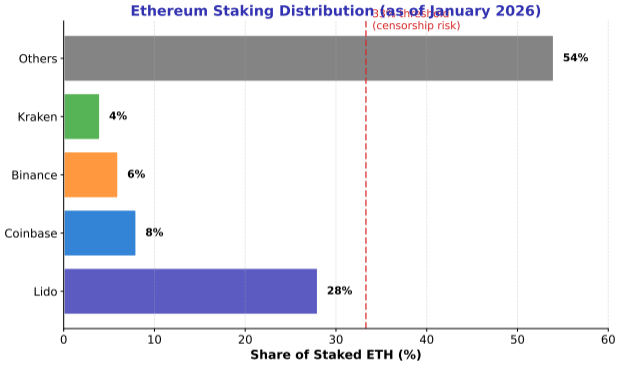
**Purpose:** Make attacks economically destructive.

---

makes the cost of attack explicit and immediate

Slash

# Staking Centralization Risks



staking protocols like Lido concentrate stake—33% threshold enables censorship

Liqui

# PoW vs PoS Comparison

## Consensus Mechanisms: PoW vs PoS

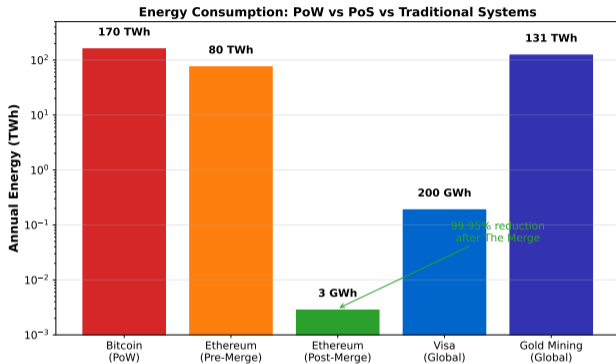
Proof of Work		Proof of Stake	
Resource:	Computing Power	Resource:	Staked Tokens
Security:	Energy Cost	Security:	Economic Loss
Hardware:	ASICs/GPUs	Hardware:	Standard Server
Attack Cost:	51% hash rate	Attack Cost:	51% stake
Example:	Bitcoin	Example:	Ethereum

VS

PoW: Spend energy to earn | PoS: Risk capital to earn

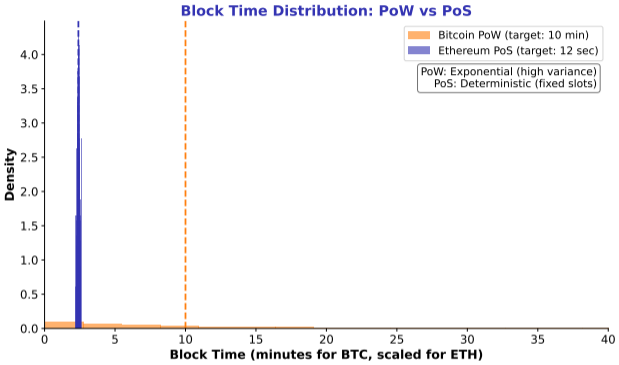
spend energy to earn, PoS: risk capital to earn

PoW:



reduced energy use by 99.995% switching from PoW to PoS

# Block Time Consistency: PoW vs PoS

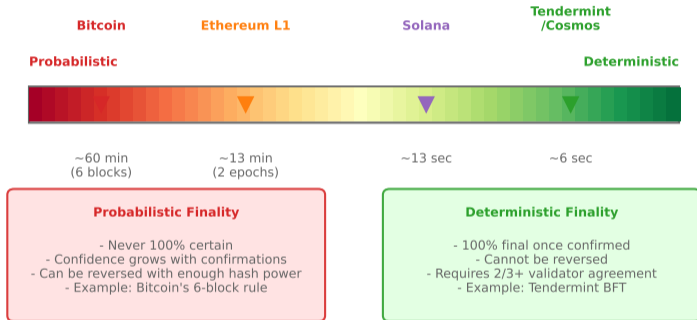


PoW

follows exponential distribution (high variance), PoS uses fixed time slots (predictable)

# Finality: When Is a Transaction Final?

## Finality Spectrum: When is a Transaction Final?

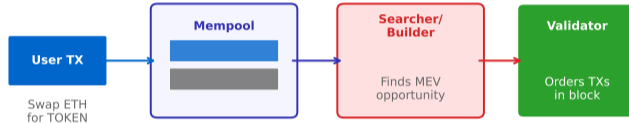


*Trade-off: Faster finality often requires more trust assumptions*

determines when you can safely consider a transaction irreversible

Finali

## Maximal Extractable Value (MEV) in PoS



### Common MEV Strategies:



MEV extracted: >\$1.5B cumulative since 2020 on Ethereum (Source: Flashbots/mevwatch.info)

is “invisible tax” on users—validators profit from transaction ordering

MEV

## What is MEV?

- Value extractable by manipulating transaction order
- Validators/builders can reorder, insert, or censor TXs
- Previously called “Miner Extractable Value”

## Impact on Users:

- Sandwich attacks: Users get worse prices
- Frontrunning: Arbitrageurs take profit first
- Failed transactions: Competing bots cause reverts

## Solutions (Ethereum):

- **Flashbots**: Private mempool, fair ordering
- **MEV-Boost**: Separates builder and proposer roles
- **Intent-based systems**: Hide transaction details

---

is a fundamental challenge—cannot be eliminated, only redistributed

MEV

## Other Mechanisms

# Delegated Proof of Stake (DPoS)

## How It Works:

- Token holders vote for delegates
- Fixed number of delegates produce blocks
- Delegates can be voted out

## Examples: EOS, Tron, Lisk Trade-offs:

- + Fast finality, high throughput
- More centralized (few delegates)
- Potential for collusion

---

sacrifices decentralization for performance

DPoS

# Practical Byzantine Fault Tolerance (PBFT)

## How It Works:

- Pre-defined set of validators
- Three-phase voting protocol
- Tolerates up to  $1/3$  malicious nodes

## Examples: Hyperledger Fabric, Tendermint Trade-offs:

- + Fast, deterministic finality
- + No energy waste
- Requires known validator set
- Does not scale well

---

is ideal for permissioned blockchains

PBFT

# Proof of Authority (PoA)

## How It Works:

- Validators are pre-approved entities
- Identity is at stake, not tokens or work
- Used in private/consortium chains

## Examples: VeChain, Ethereum testnets Trade-offs:

- + Very fast and efficient
- + Low hardware requirements
- Centralized trust
- Not suitable for public chains

---

prioritizes efficiency over decentralization

PoA

## Choosing a Consensus Mechanism

## Cannot optimize all three simultaneously:

- ① **Security:** Resistance to attacks
- ② **Decentralization:** No single point of control
- ③ **Scalability:** High transaction throughput

## Trade-offs by Mechanism:

- PoW: Security + Decentralization (sacrifices scalability)
- DPoS: Security + Scalability (sacrifices decentralization)
- PBFT: Security + Scalability (sacrifices decentralization)

---

consensus mechanism makes trade-offs

Every

## What You Learned Today:

- 1 Byzantine Generals Problem: Agreement with potential traitors
- 2 Proof of Work: Spend energy for security
- 3 Proof of Stake: Risk capital for security
- 4 Trade-offs: Security vs. decentralization vs. scalability

**Core Insight:** Consensus mechanisms convert some scarce resource (energy, capital, reputation) into network security.

## Questions for Reflection

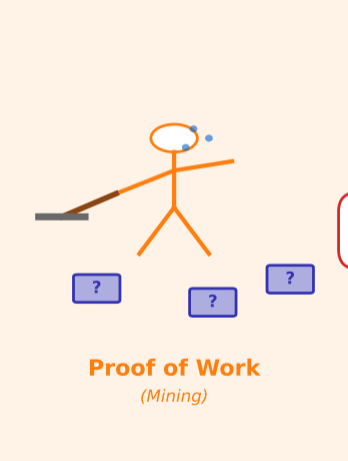
- 1 Why is PoW called “one CPU one vote”?
- 2 How does slashing make PoS secure?
- 3 When would you choose DPoS over PoW?
- 4 Can the blockchain trilemma be solved?

**Discussion:** Is PoW's energy consumption justified for security?

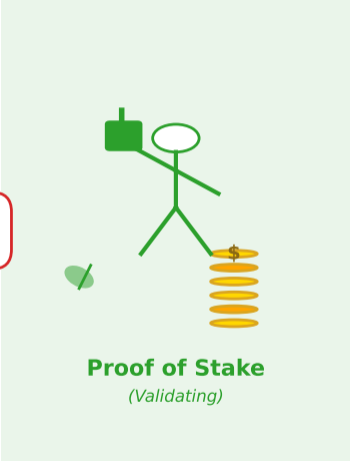
---

these questions before our next session

Consi



VS



## Primary References for Course Data:

- 1 Cambridge Bitcoin Electricity Consumption Index (CBECI)
- 2 Ethereum Foundation: PoS Documentation
- 3 Digiconomist: Bitcoin Energy Consumption Tracker
- 4 Blockchain.com: Bitcoin Network Hash Rate Charts
- 5 MEV Watch: Real-time MEV Monitoring

All data current as of January 2026. Energy estimates include both direct mining and cooling infrastructure.

---

data sources independently for research purposes

Verify

Thank You

Questions?

Course materials: [digital-ai-finance.github.io/crypto-economics](https://digital-ai-finance.github.io/crypto-economics)

## Appendix: Technical Deep Dives

# Case Study: The Ethereum Merge (September 2022)

## What Happened:

- Ethereum switched from PoW to PoS
- Largest blockchain upgrade in history
- Zero downtime during transition

## Results:

- Energy reduction: 99.95% (Ethereum Foundation estimate) (from 100 TWh to 0.01 TWh/year)
- New ETH issuance dropped 90%
- Created world's largest staking economy

## Technical Achievement:

- Beacon Chain ran parallel since Dec 2020
- Difficulty bomb forced transition
- No “replay” of transactions needed

---

Merge proved major PoW chains can transition to PoS safely

The

### The Problem:

In PoS, validators can vote on multiple chains for free (no work required).

### Why It's Dangerous:

- Validators can hedge by supporting all forks
- No cost to supporting conflicting histories
- Could prevent consensus from forming

### Solutions:

- **Slashing:** Destroy stake for double-voting (Ethereum)
- **Deposit lockup:** Can't withdraw during dispute window
- **Checkpointing:** Finalize blocks periodically

### Ethereum's Approach:

Validators who attest to conflicting blocks lose 1/32 of stake initially, up to 100% if attack is widespread.

---

converts "nothing at stake" to "everything at stake"

Slash

### The Attack:

- Attacker acquires old private keys (from early validators)
- Creates alternative history from early block
- New nodes can't distinguish real vs fake chain

### Why PoS is Vulnerable:

- No work to rewrite history (unlike PoW)
- Old keys have no current stake (already withdrawn)
- Can create plausible-looking fork

### Solutions:

- **Weak subjectivity:** Trust recent checkpoint
- **Key deletion:** Validators must delete old keys
- **Social consensus:** Community rejects obvious attacks

---

range attacks are theoretical—require social engineering to succeed

Long-

## Technical: Consensus Mechanisms Comparison

Property	PoW	PoS	DPoS	PBFT
Validators	Anyone	Stakers	Elected	Permissioned
Finality	Probabilistic	Probabilistic*	Fast	Immediate
TPS	3-7	15-30	1000+	1000+
Energy	Very High	Low	Low	Low
51% Attack	Hash power	Stake	Stake+Votes	1/3 nodes
Decentralization	High	Medium-High	Medium	Low
Example	Bitcoin	Ethereum	EOS	Hyperledger

\*Ethereum PoS provides economic finality via slashing after 2 epochs ( 13 min)

perfect consensus mechanism exists—each makes different trade-offs

No