

Catastrophe by Design

Destabilizing Wasteful Technologies & The Phase Transition from
Proof of Work to Proof of Stake

Stefanos Leonardos, Iosif Sakos, Costas Courcoubetis and
Georgios Piliouras



Energy Consumption of PoW

Proof of Work (PoW): Security \iff Work \iff Energy Consumption

- 1 BTC transaction = 775.818 VISA transactions.
- BTC consumes more energy than Finland and Pakistan.
- Energy consumption **doubles** every year.
- BTC is only one out of many PoW blockchains, e.g., **Ethereum**.

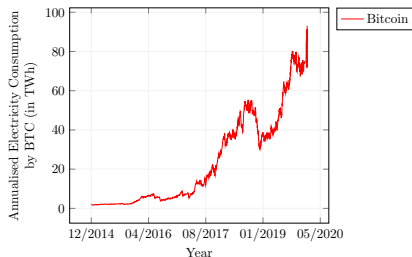


Figure: The Cambridge Bitcoin Electricity Consumption Index (CEBCI)

Transition to PoS

Proof of Stake (PoS) has equivalent provable guarantees to PoW. But:

- More **work** implies more **safety** – more reliable applications (e.g., BTC).
- When all PoW it is **individually rational** to also PoW.
- Even worse PoW is evolutionary stable: **small groups of adopters** of alternative technologies are doomed to fail.

These observations hint towards a game-theoretic model.

Model I: Agents and Strategies

A **population** p of agents, investors or **miners** (physical or virtual)

- Mass $K > 0$: total available **capital or resources**, e.g., money, hardware or electricity.
- **Strategies**: two available technologies, W (costly), and S .
- Investment cost: $\gamma > 0$ for W and 0 for S .
- **Population states**: $X = \{(x, 1 - x) : x \in [0, 1]\}$ where x = fraction of **PoW investors**

Model II: Value and Payoffs

Each technology creates **value** split among **adopters**

- **Value** V , **Adoption** $\alpha > 1$:
 - $V_W = V(xK)^\alpha$ and $V_S = V((1-x)K)^\alpha$
- **Payoff functions**: equal share amongst all invested units:
 - $u(W, x) = V_W \cdot (xK)^{-1} - \gamma = VK^{\alpha-1}x^{\alpha-1} - \gamma$
 - $u(S, x) = V_S \cdot ((1-x)K)^{-1} = VK^{\alpha-1}(1-x)^{\alpha-1}$
- For the purposes of this talk we restrict ourselves to the case $\alpha = 2$:
 - $u(W, x) = VKx - \gamma$
 - $u(S, x) = VK(1-x)$

An Evolutionary Game

Evolutionary game interpretation

$$P = \begin{array}{c} W \\ S \end{array} \begin{array}{cc} W & S \\ \left(\begin{array}{cc} VK - \gamma & -\gamma \\ 0 & VK \end{array} \right) \end{array} \quad (G1)$$

Theorem

*(G1) has three **Nash equilibria**: (W, W) , (S, S) and one mixed. The two pure equilibria are **evolutionary stable**, whereas the mixed one is **unstable**.*

Population Dynamics

Q-Learning dynamics:

$$\dot{x} = x \left[\underbrace{u(W, x) - \bar{u}(x)}_{\text{Replicator Dynamics}} - T \cdot \underbrace{(x \ln x + (1 - x) \ln (1 - x))}_{\text{Entropy}} \right]$$

Where $\bar{u}(x) = xu(W, x) + (1 - x)u(S, x)$

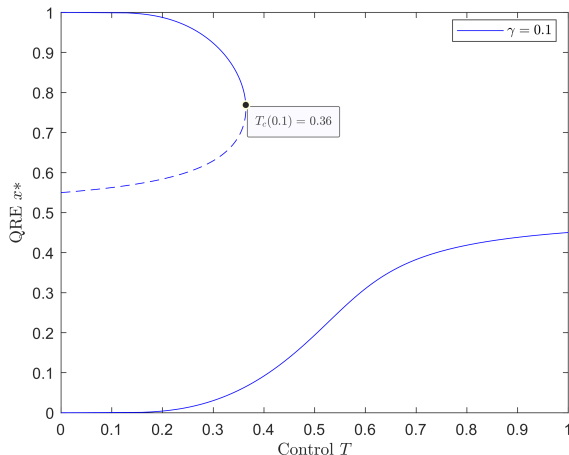
Quantal Response Equilibrium (QRE): The **steady** states of the system, i.e., $\dot{x} = 0$.

We can affect the agents' **rationality** by **scaling** the agents utilities:

$$\begin{aligned} x \left[\frac{u(W, x)}{c} - \frac{\bar{u}(x)}{c} - T \cdot (x \ln x + (1 - x) \ln (1 - x)) \right] &= 0 \\ \iff x [u(W, x) - \bar{u}(x) - cT \cdot (x \ln x + (1 - x) \ln (1 - x))] &= 0 \end{aligned}$$

QRE Correspondence: Visually

In our case ($\alpha = 2$): $\dot{x} = x(1-x)[2x - (1 + \gamma) - T \ln(\frac{x}{1-x})]$



QRE Correspondence: Formally

Theorem

For any $\alpha > 1$ there exists a finite sequence of temperatures $T = \langle T_0, T_1, \dots \rangle$ such as starting from an initial state x_0 and performing the following procedure for each $T_i \in T$:

- *Scale the system's temperature at T_i , and*
- *Wait until the system converges to a QRE*

the system is going to converge to the desirable state $x = 0$ which corresponds to energy-friendly technology S .

We can reliably **destabilize PoW** equilibrium and **converge to PoS** equilibrium by introducing and removing **taxes** in the system.

Short Term Policy \implies Long Lasting Effects

Conclusion

