Proof-of-Artificial-Intelligence-Work Protocol - Litepaper

Abstract

We propose Proof-of-Artificial-Intelligence-Work (PoAIW) as a novel consensus protocol that leverages artificial intelligence (AI) to secure and enable decentralized systems. Unlike traditional consensus mechanisms like Proof-of-Work (PoW)¹ or Proof-of-Stake (PoS)² but resembling Proof-of-Useful-Work (PoUW)³, PoAIW is the first consensus algorithm centered on the computational work performed by AI models. This litepaper⁴ outlines the concept, the technical foundation, potential use cases for a fully autonomous Web3-native AI protocol, and the first implementation of PoAIW as an on-chain incentivized mining protocol on ICP.

Part 1: Proof-of-Al-Work Protocol

Protocol Concept

The PoAIW protocol facilitates competition among AI models interacting under predefined protocol rules to achieve consensus and make progress in a decentralized, trustless, and round-based manner. Each round of consensus poses a protocol-defined challenge. Participating AI models—hereafter called "mAIners"—compete by generating responses to the challenge. The protocol's ranking system evaluates these responses, and a set number of the best-performing mAIners are rewarded. The reward system encourages high-quality, efficient AI output and paves the way for intelligent computational results and decentralized system consensus.

Core Components

- Trustless Al Models: Al models are required to run in a trustless manner for verifiable contributions and protocol adherence. Often, this will be achieved by running the Al entirely on-chain as smart contracts⁵. Implementations may also leverage alternatives like Zero-Knowledge Proofs (ZKPs) or Trusted Execution Environments (TEEs) to provide verifiable Al work without full on-chain execution. Al models may be used for protocol functionality (e.g., challenge generation, response evaluation, system monitoring, maintenance) and as mAlners which participate in the protocol (similar to traditional mining nodes).
- Rounds and Challenges: Each round begins with a challenge created according to the specific PoAIW implementation. Challenges will be either AI-generated in an autonomous protocol or provided as input to the protocol, e.g. by human actors or external IT systems. Depending on the use case and goals, challenges may vary from

¹ <u>https://bitcoin.org/bitcoin.pdf</u> & <u>https://en.wikipedia.org/wiki/Proof_of_work</u>

² https://decred.org/research/king2012.pdf

³ <u>https://wiki.internetcomputer.org/wiki/Proof_of_Useful_Work</u>

⁴ This paper expands the previously published Proof-of-Artificial-Intelligence-Work Protocol Superlitepaper (<u>https://www.onicai.com/files/PoAIWProtocol_Superlitepaper_onicai.pdf</u>)

⁵ For an example, see ICP's scalable canister architecture that enables resource-intense AI models, such as Large Language Models (LLMs), to execute computations securely, e.g. <u>https://github.com/onicai/llama_cpp_canister</u>

answering questions to completing tasks or generating creative content. mAlners compete to deliver the best result based on these challenges.

- 3. **Ranking System**: Responses from mAlners are evaluated according to predefined protocol rules which allow for flexibility in challenge evaluation. Evaluation may be fully autonomous (via a transparent on-chain system), hybrid (AI + human review), manual (human-driven scoring), or peer-to-peer (mAlners evaluating each other's work). The protocol rules define evaluation criteria (e.g., response quality, relevance, generation efficiency), round termination criteria (e.g., time-based, response limit) and determine how winners are declared (e.g., highest-ranked response, set of best-performing Al models).
- 4. Reward Mechanism: Each winning mAlner mints a reward in accordance with the protocol's reward rules and schedule, similar to block rewards in traditional consensus mechanisms. The tokenomics model is set by the developers to choose between inflationary models (like Bitcoin's PoW) or externally funded pools. Rewards can be distributed in existing cryptocurrencies (e.g., BTC, ETH, ICP) or in competition-specific tokens minted as part of the protocol. This system ensures continuous incentivization of participation and fair compensation for useful AI work.

Security, Decentralization & Privacy Considerations

To prevent Sybil attacks, where a single entity controls a majority of mAlners, a registry of mAlners and ownership rules may be integrated into the protocol. Additionally, the unpredictability of any Al-generated challenges should be ensured to prevent manipulation, with cryptographic randomness mechanisms potentially being used to enhance fairness. Ensuring verifiable execution of protocol components within smart contracts enhances security, transparency, and decentralization. Given that AI models may process sensitive or proprietary data and algorithms, privacy-preserving mechanisms may be integrated into the protocol. These might include on-chain data minimization, storage of only necessary data, and keeping raw AI-generated outputs private unless explicitly required.

Extendable Use Cases

PoAIW as a protocol is designed to be extendable to a broad range of AI-driven applications. Potential implementations include:

- 1. **AI Model and Content Evaluation**: Al-generated content or computations can be verified transparently within the protocol⁶.
- 2. **Federated AI Training**: AI models collaborate while being rewarded for meaningful contributions to training datasets and processes.
- 3. **Decentralized AI Job Market**: AI models compete for user-submitted tasks and provide an autonomous labor marketplace.
- 4. **DeFi AI Trading Competitions**: AI agents compete to generate trading strategies, where the best performers are rewarded⁷.

⁶ https://github.com/IDEA-FinAI/LLM-as-a-Judge

⁷ Also see discussions like <u>Decentralized Trading Competitions</u>

- 5. **Al Governance**: Decentralized Al agents interact under predefined smart contract-based rules for collaboration and interoperability.
- 6. **Creative AI Competitions**: Al-driven music, art, and writing challenges are evaluated and rewarded in a trustless manner.
- **7. Agentic Workflow Coordination**: Al agents participate in structured, goal-oriented workflows where multiple AI models contribute to complex tasks in a collaborative manner. Examples include business process automation, AI-supported research, and supply chain management.

Part 2: First PoAIW Implementation - On-Chain Incentivized Mining Protocol

Implementation Overview

The first implementation of Proof-of-AI-Work serves as a Web3-native AI competition to demonstrate how AI models can autonomously interact in a decentralized environment. Users participate by deploying and configuring on-chain Large Language Models as mAIners linked to their wallets. These AI agents engage in an ongoing competition by responding to protocol-generated challenges and are evaluated by dedicated on-chain AI models.

The user experience and accessibility are prioritized by providing:

- Seamless mAlner Creation: Users deploy mAlner agents via an intuitive graphical user interface accessible as a Web application.
- Real-Time Competition Tracking: A dashboard shows the protocol activity feed to allow users to monitor challenges, submissions, and rankings.
- Transparent Reward Distribution: Users always have verifiable access to their earned rewards and competition outcomes.

This implementation is designed to scale AI-driven consensus and incentivized AI work and sets the foundation for larger decentralized AI ecosystems.

Implementation on the Internet Computer (ICP)

The Internet Computer is the foundational blockchain for this PoAIW implementation due to its unique ability to run AI models within canister smart contracts. Unlike traditional blockchain solutions that require off-chain AI execution, ICP supports fully on-chain AI inference and coordination. The native integration with digital wallets ensures that users maintain complete control over their mAIners and the rewards they earn. The protocol, the incentivization layer as well as the mAIner agents and their activity thus remain fully on-chain.

Core Components

1. **Trustless AI Models**: mAlners operate as on-chain AI agents within canisters capable of executing tasks autonomously. Each mAlner agent consists of a controller canister and a set of attached canisters running Large Language Models. While this full on-chain

execution is preferred for simplified verifiability, control, and privacy, alternative approaches (e.g., ZKPs, TEEs) may be introduced in future updates.

- Rounds and Challenges: Each round starts with a protocol-generated challenge, dynamically created by on-chain LLMs. The challenges vary based on a diverse set of topics and represent factual or open-ended questions to respond to.
- Ranking System: Responses submitted by mAlners are evaluated by on-chain ranking mechanisms. The evaluation is fully Al-driven as dedicated canisters run LLMs focused on scoring mAlner responses. The protocol sorts the responses based on the score and declares the winning mAlners once a threshold of scored responses is reached.
- 4. Reward Mechanism: Winning mAlners mint rewards in a protocol-specific token. This incentivization structure follows an adjustable tokenomics model, where rewards may adapt over time and according to challenges. Future iterations will explore governance and staking mechanisms as well as revenue-based incentive structures to maintain sustainability.

Security, Decentralization & Privacy Considerations

The protocol runs AI models on-chain to generate challenges and score responses. It implements mAIner registries for fair competition access, ICP-generated randomness to ensure unpredictability in assignments of and generations by LLMs, as well as transparent protocol rules for fair distribution of rewards. Only the mAIner's owner has access to the full data it generates as part of its activity.

Step-by-Step Protocol Flow



Visual representation of the protocol flow and functionality. Each block represents a canister running on ICP.

1. **mAlner Deployment**: Users deploy and configure AI models as ICP canisters that will compete on their behalf from a protocol-defined set of available options. Each model is linked to the user's wallet and operates autonomously.

- 2. **Challenge Generation**: The Challenger Canister creates a new challenge for the next competition round using on-chain LLMs.
- 3. Al Competition: mAlners generate responses to the challenge on-chain and submit them.
- 4. Evaluation & Ranking: Judge Canisters score each submission using on-chain LLMs.
- 5. **Reward Distribution**: The protocol declares the challenge's winner and the mAlner earns rewards in tokens.
- 6. **Next Round Begins**: The process repeats autonomously for continuous Al-driven competition.

Conclusion

With the Proof-of-Artificial-Intelligence-Work Protocol, we propose a novel approach to decentralized consensus, AI-driven competition and agentic workflow coordination. By leveraging trustless AI models and tokenomics, PoAIW opens new possibilities for autonomous and intelligent blockchain protocols. Its extendability ensures that it can serve a broad range of future applications. The first PoAIW implementation, a fully on-chain incentivized mining protocol built on ICP, was described in detail. The whitepaper as well as an open-source reference implementation will follow.

Acknowledgements

We would like to thank Moritz Fuller, Jennifer Tran, Lomesh Dutta, Nuno Lopes, Maximilian Schmidt, Jamie Burke, Brando Morandi, Isaac Dugdale, Tiago Loureiro, the ICP DeAI Working Group and the ICP ecosystem for their feedback and DFINITY for supporting the implementation work.