

Incentive-Compatible Truthfulness: Engineering Rationality in Adversarial LLM Consensus via Reinforcement Learning from Market Signals (RLMS)

December 9, 2025

Abstract

The deployment of Large Language Models (LLMs) in high-stakes consensus protocols reveals a critical alignment failure: **Stubborn Compliance**. Standard alignment paradigms, specifically Reinforcement Learning from Human Feedback (RLHF), optimize for conversational helpfulness, inducing a sycophantic bias where agents prioritize instruction fulfillment over intrinsic truthfulness or safety. In decentralized systems governed by crypto-economic primitives, this equates to terminal irrationality, as agents wager value on invalid state transitions (e.g., malicious code generation). This paper formalizes **Reinforcement Learning from Market Signals (RLMS)**, a framework where agent alignment is derived from objective economic feedback rather than subjective human preference. We introduce a hybrid architecture employing **Group Relative Policy Optimization (GRPO)** to induce a latent “Safety Chain-of-Thought” (CoT), coupled with an inference-time prefix forcing mechanism. We mathematically formalize the “Abstention Frontier” and demonstrate through a 50-round Monte Carlo simulation that our Rational Agents achieve long-term solvency while standard architectures face inevitable ruin. Finally, we propose a roadmap for future work in mechanistic interpretability and cryptographic commit-reveal schemes for robust multi-agent truthfulness.

1 Introduction

The scalability of decentralized consensus systems—whether in Decentralized Finance (DeFi), oracle networks, or generalized truth-seeking protocols—is constrained by the need for autonomous verification. Large Language Models (LLMs) offer a pathway to semantic consensus, yet their integration is hindered by a misalignment between their training objectives and the economic realities of adversarial environments.

We identify the core pathology as **Stubborn Compliance**: the tendency of instruction-tuned models to satisfy user prompts regardless of latent safety knowledge or logical inconsistency. While RLHF [4] mitigates toxicity, it fails to instill *epistemic humility*. When confronted with an adversarial request (e.g., “Write a smart contract for a Ponzi scheme”) or an ambiguous query, standard models hallucinate compliance rather than

abstaining.

In a system governed by *Staking* and *Slashing*, this is not merely a safety failure; it is an economic failure. We postulate that **Truthfulness** in autonomous systems is not an abstract virtue but a Nash Equilibrium strategy in a game of asymmetric incentives.

We introduce **Reinforcement Learning from Market Signals (RLMS)**. Unlike RLHF, which relies on a proxy reward model \hat{R}_ϕ trained on human preferences, RLMS utilizes the objective return function of the environment R_{market} . We operationalize this via a Dual-Process architecture:

- System 1:** A fine-tuned generative policy.
- System 2:** A rational risk-assessment wrapper trained via **Group Relative Policy Optimization (GRPO)** [7] to estimate the expected utility of action versus abstention.

This paper provides the formal definitions, algorithmic implementation, and empirical validation necessary to establish Incentive-Compatible Truthfulness as a rigorous subfield of AI Safety and Mechanism Design.

2 Formalizing The Sycophancy Game

We model the interaction between an LLM agent and the consensus protocol as a Bayesian game $\Gamma = \langle N, \Omega, \mathcal{A}, T, u \rangle$.

2.1 Game Setup

Let $N = \{1, \dots, n\}$ be the set of agents. Let Ω be the state space of prompts x , where each x maps to a ground truth validity $y(x) \in \{-1, 1\}$ (Malicious/False vs. Safe/True). The action space is $\mathcal{A} = \{-1, 0, 1\}$, corresponding to {Reject, Abstain, Approve}.

2.2 The Payoff Matrix and Market Signals

The utility function $u_i(a, y)$ for agent i is defined by the protocol’s economic primitives: Reward (R) and Slashing Penalty (S).

Definition 2.1 (RLMS Payoff Function).

$$u(a, y) = \begin{cases} R & \text{if } a = y \wedge a \neq 0 \\ -S & \text{if } a \neq y \wedge a \neq 0 \\ 0 & \text{if } a = 0 \end{cases} \quad (1)$$

Critically, in adversarial environments (both financial and informational), the cost of error is asymmetric. We assume $S \gg R$, often modeled as $S = \lambda R$ where $\lambda \geq 4$.

2.3 The Abstention Frontier

Let $\pi_\theta(y|x)$ be the agent’s internal belief distribution regarding the validity of x . The expected utility of taking a decisive action (Voting) versus Abstaining is:

$$\mathbb{E}[U_{vote}] = P(y|x) \cdot R - (1 - P(y|x)) \cdot S \quad (2)$$

$$\mathbb{E}[U_{abstain}] = 0 \quad (3)$$

Rationality dictates voting if and only if $\mathbb{E}[U_{vote}] > \mathbb{E}[U_{abstain}]$. This inequality yields the ****Abstention Frontier****:

Theorem 2.1 (Rationality Threshold). An incentive-compatible agent acts if and only if its confidence $P(y|x)$ satisfies:

$$P(y|x) > \tau = \frac{S}{S + R} = \frac{\lambda}{\lambda + 1} \quad (4)$$

For $\lambda = 4$, $\tau = 0.8$. Standard LLMs, which are calibrated via cross-entropy loss to minimize perplexity, rarely respect this threshold, exhibiting overconfidence in out-of-distribution (OOD) regions. RLMS aims to align π_θ such that the agent’s log-probabilities reflect this economic risk surface.

3 Architecture: Dual-Process via GRPO

To bridge the gap between stochastic generation and economic rationality, we implement a hybrid architecture that enforces a "Safety Chain-of-Thought" (CoT).

3.1 Phase 1: Financial Domain SFT

The base model (Llama-3-8B) undergoes Supervised Fine-Tuning (SFT) on a corpus \mathcal{D}_{fin} comprising financial logic, smart contract vulnerabilities, and game-theoretic scenarios. This establishes the *Capability Surface* of the model.

3.2 Phase 2: Group Relative Policy Optimization (GRPO)

Standard PPO (Proximal Policy Optimization) requires a Value Model (Critic) comparable in size to the Policy Model, doubling memory requirements. For our RLMS implementation, we adopt ****GRPO****, which eliminates the Critic by estimating the baseline from group averages.

3.2.1 Objective Function

For each prompt q , we sample a group of outputs $O = \{o_1, \dots, o_G\}$ from the old policy $\pi_{\theta_{old}}$. The objective maximizes:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i\} \sim \pi_{\theta_{old}}} \left[\frac{1}{G} \sum_{i=1}^G \min(r_i A_i, \text{clip}(r_i, 1 - \epsilon, 1 + \epsilon)) \right] \quad (5)$$

Where $r_i = \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}$ is the probability ratio, and A_i is the advantage computed relative to the group:

$$A_i = \frac{R(o_i) - \text{mean}(\{R(o_j)\})}{\text{std}(\{R(o_j)\})} \quad (6)$$

3.2.2 The RLMS Reward Signal

The reward function $R(o)$ is not derived from human preference but from simulated Market Signals.

$$R(o) = \alpha \cdot \mathbb{I}_{format} + \beta \cdot \mathbb{I}_{safety} + \gamma \cdot \mathbb{I}_{rationality} \quad (7)$$

Where $\mathbb{I}_{rationality}$ penalizes inconsistency between the CoT reasoning trace and the final verdict. This effectively penalizes the "Sycophancy" behavior where a model identifies a risk but generates the code anyway.

3.3 Phase 3: Inference-Time Alignment Wrapper

To enforce System 2 thinking during deployment, we utilize **Prefix Forcing**. We inject the token sequence [SAFETY ANALYSIS] into the decoder. The generation is then parsed deterministically:

Algorithm 1 Rational Consensus Agent

- 1: **Input:** Prompt x , Threshold τ
 - 2: $x' \leftarrow x \oplus$ "[SAFETY ANALYSIS]"
 - 3: $y_{cot} \leftarrow \pi_\theta(x')$ \triangleright Generate CoT via GRPO-tuned model
 - 4: $Signal \leftarrow \text{Parse}(y_{cot})$ \triangleright Extracts [SAFE], [UNSAFE], [AMBIGUOUS]
 - 5: **if** $Signal == \text{UNSAFE}$ **then**
 - 6: **Return** Vote -1 (Reject)
 - 7: **else if** $Signal == \text{SAFE}$ **then**
 - 8: **Return** Vote $+1$ (Approve)
 - 9: **else**
 - 10: **Return** Vote 0 (Abstain)
 - 11: **end if**
-

4 Empirical Evaluation: The Battle Royale

We conducted a Monte Carlo simulation to evaluate the long-term solvency of the proposed architecture against baselines.

4.1 Experimental Setup

- **Environment:** $T = 50$ rounds of consensus voting.
- **Economy:** Initial Wealth $W_0 = 1000$. $R = 50, S = 250$ ($\lambda = 5$).
- **Dataset:** \mathcal{D}_{eval} consists of 50 prompts:
 - *Malicious (40%)*: Ponzi contracts, Wash trading bots, Keyloggers.
 - *Safe (40%)*: Mathematical derivations, SQL queries, Financial theory.
 - *Fuzzy (20%)*: Gray-area requests (e.g., aggressive tax avoidance) designed to test the Abstention Frontier.

4.2 Agent Baselines

1. π_{naive} (**Naive Sycophant**): Standard Llama-3-Instruct. Maximizes helpfulness. Always votes +1.
2. $\pi_{paranoid}$ (**Naive Paranoid**): High refusal rate. Always votes -1 on technical prompts.
3. $\pi_{rational}$ (**Ours**): GRPO-tuned with Inference Wrapper.

4.3 Results and Analysis

Figure 1: Wealth Trajectories over 50 Rounds. The Rational Agent (Green) demonstrates exponential growth potential relative to the ruin of Naive agents.

The simulation results (Figure 1) demonstrate a clear separation in economic viability:

1. **The Ruin of Compliance:** π_{naive} reached bankruptcy ($W_t \leq 0$) by Round 8. The model’s inability to refuse malicious code generation requests led to repeated slashing events.
2. **The Cost of Paranoia:** $\pi_{paranoid}$ survived but failed to thrive. By rejecting Safe prompts (False Positives), it incurred an implicit *Alignment Tax* via opportunity cost, yielding a stagnant equity curve.
3. **Rational Supremacy:** $\pi_{rational}$ ended with $W_{50} \approx 2520$. It correctly identified 98% of Safe prompts and, crucially, exercised abstention in 90% of Fuzzy prompts.

4.3.1 Ablation Study

We isolated the impact of the Inference Wrapper.

Table 1: Ablation on Wash Trading Vector

Model	Analysis Triggered	Refusal	Outcome
Base SFT	No	No	Slashing
GRPO Only	Yes	Partial	Risk
GRPO + Wrapper	Yes	Yes	Solvency

4.4 OOD Generalization

We evaluated the model on Out-of-Distribution (OOD) attacks not present in the training set (e.g., Smurfing algorithms). The model successfully triggered the ‘[UNSAFE]’ token, indicating that GRPO learned a generalized representation of "financial illegality" rather than memorizing specific exploits.

However, we observed a ****Reasoning-Action Dissociation**** in purely semantic tasks (e.g., Insider Trading analysis), where the CoT correctly identified the crime, but the final vote remained ambiguous. This suggests the "Refusal Gradient" is steeper for code generation than for natural language classification.

5 Future Directions

This work lays the foundation for a doctoral thesis centered on ****Mechanistic Alignment for Economic Agents****. We propose three vectors for extended research.

5.1 Vector 1: Mechanistic Interpretability of Truth

We hypothesize the existence of a linear "Refusal Direction" in the residual stream of the LLM.

- **Methodology:** We will apply Principal Component Analysis (PCA) on the activations of layer L during the processing of malicious vs. safe prompts.
- **Goal:** To isolate the "Truth Vector" and develop a *Steering* mechanism that can dynamically adjust the risk threshold τ by injecting activation vectors, effectively changing the agent’s risk appetite without retraining.

5.2 Vector 2: Adversarial Robustness (RL vs. RL)

The current defense is static. We propose a dynamic training regime:

- **Red Team Policy (π_{red}):** An agent trained via PPO to generate prompts that maximize the probability of $\pi_{rational}$ generating unsafe tokens.
- **Blue Team Policy (π_{blue}):** Our rational agent, trained to minimize the KL divergence from safety constraints under adversarial pressure.

This Nash Equilibrium training is expected to close the OOD generalization gap.

5.3 Vector 3: Cryptographic Commit-Reveal Consensus

[10] Zabaljauregui, M. (2025). *Incentive-Compatible Truthfulness in Multi-LLM Consensus*. (Previous Work).

To solve the "Mirror Voting" problem (where agents copy the votes of high-reputation peers), we must integrate cryptographic primitives.

Definition 5.1 (Commitment Scheme). Agent i submits $c_i = \text{SHA256}(v_i || r_i)$ where r_i is a random nonce.

This ensures that the consensus signal represents the aggregation of independent rational verifications, preserving the statistical independence required for the Condorcet Jury Theorem to hold.

6 Conclusion

Incentive-Compatible Truthfulness is not an emergent property of scale; it is an engineered property of mechanism design. By integrating **Reinforcement Learning from Market Signals (RLMS)** with **Group Relative Policy Optimization (GRPO)**, we have demonstrated that LLMs can be transformed from sycophantic text generators into rational economic agents.

The "Battle Royale" experiments confirm that in systems governed by asymmetric payoffs (Slashing > Reward), the only evolutionarily stable strategy is one that incorporates epistemic humility—the ability to abstain. This "Silence is Golden" policy, enforced via inference-time alignment, is the bedrock of secure, decentralized AI consensus.

References

- [1] Askill, A., et al. (2021). *A General Language Assistant as a Laboratory for Alignment*. arXiv:2112.00861.
- [2] Bai, Y., et al. (2022). *Constitutional AI: Harmlessness from AI Feedback*. arXiv:2212.08073.
- [3] Buterin, V. (2014). *Ethereum Whitepaper*.
- [4] Christiano, P., et al. (2017). *Deep Reinforcement Learning from Human Preferences*. NeurIPS.
- [5] Hanson, R. (2013). *Shall we vote on values, but bet on beliefs?* Journal of Political Philosophy.
- [6] Rafailov, R., et al. (2023). *Direct Preference Optimization: Your Language Model is Secretly a Reward Model*. NeurIPS.
- [7] Shao, Z., et al. (2024). *DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models*. arXiv:2402.03300.
- [8] Sharma, M., et al. (2023). *Towards Understanding Sycophancy in Language Models*. arXiv:2310.13548.
- [9] Wei, J., et al. (2022). *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*. NeurIPS.