

Training Language Models for Social Deduction with iliad **Multi-Agent Reinforcement Learning**

Bidipta Sarkar, Warren Xia, C. Karen Liu, Dorsa Sadigh



$$L_{\rm L}(\pi, \tau_t^i) = -\log \pi(a_{\text{vote},j} | \tau_t^i)$$

$$B_t = \sum_{k \in C_t} \pi^k (a_{\text{vote},j} | \tau_t^k)$$

$$r_t^s = B_t - B_{t'}$$



Ablation Study Basic RL improves crewmate win rate, but voting and discussing remains challenging π_{RWKV} π_{RWKV7B} Π_l π_{RL+L} π_{RL+L+S} Models

Our additional components significantly improve capabilities, nearly doubling win rates

Imposter Robustness

Iterated self-play converges, indicating that learned conventions are hard to exploit



Code & Models

socialdeductionllm.github.io

