



中國人民大學

RENMIN UNIVERSITY OF CHINA



高瓴人工智能學院

Gaoling School of Artificial Intelligence

# Reflective Multi-Agent Collaboration based on Large Language Models



# Reflective Multi-Agent Collaboration based on Large Language Models

In this paper, we consider leveraging the self-reflection mechanism to improve multi-agent collaboration, and propose an elegant framework named COPPER.

- **Agent Decision Process**

$$a_{k,t} = \text{Actor}^i(p^i, sm_{k,t}^i, s_{k,t}).$$

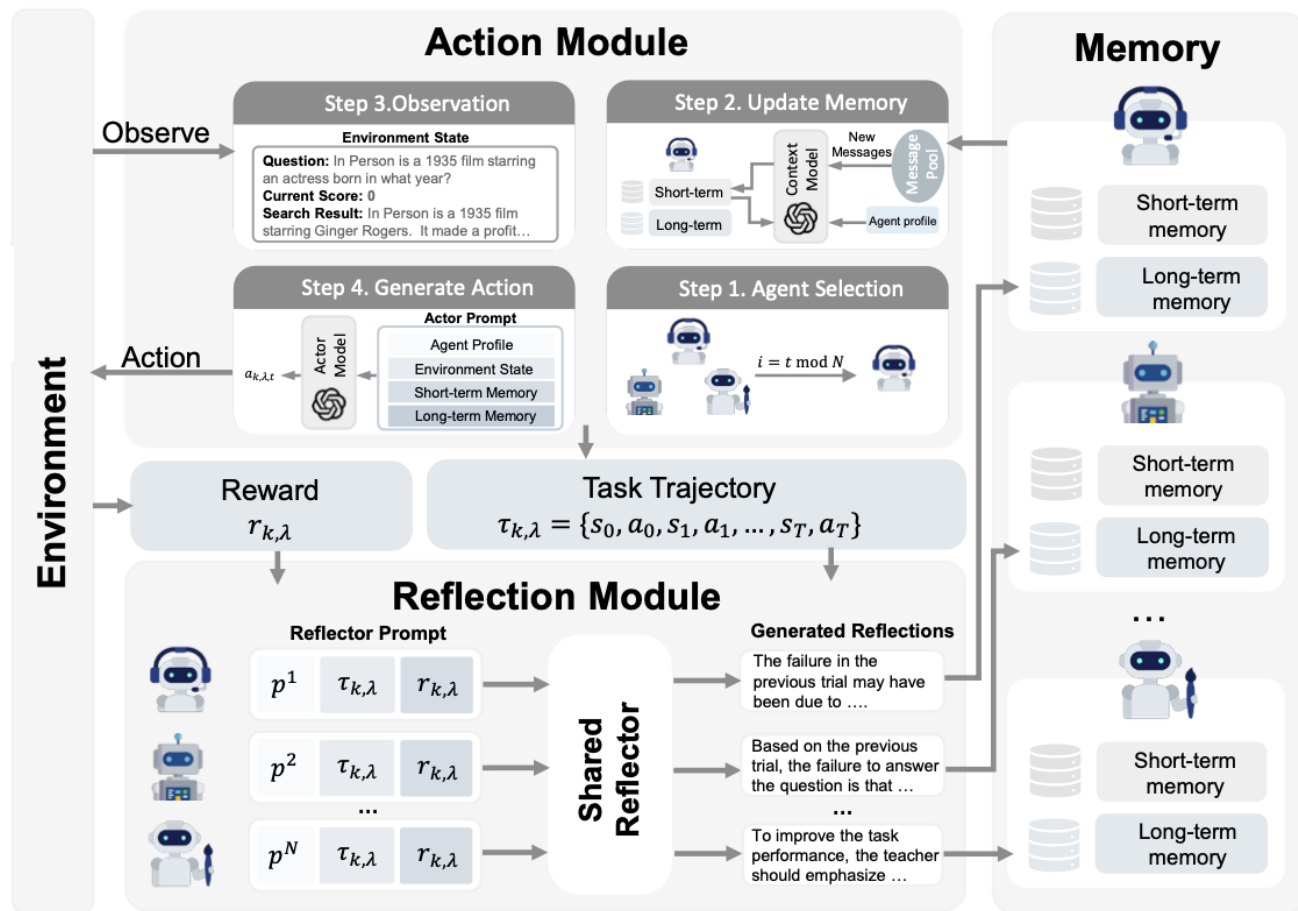
$$a_{k,\lambda,t} = \text{Actor}^i(p^i, lm_{k,\lambda}^i, sm_{k,\lambda,t}^i, s_{k,\lambda,t}),$$

- **Agent Reflection Process**

$$y_{k,\lambda}^i = \text{Reflector}^i(p^i, [sm_{k,\lambda,T}^i]_{i=1}^N, r_{k,\lambda}),$$

- **Update of Short-Term Memory**

$$sm_{k,t}^i = \text{Context}^i(p^i, sm_{k,t-1}^i, \{s_i, a_i\}_{i=\max(0,t-N+1)}^t),$$



# Reflective Multi-Agent Collaboration based on Large Language Models

- Towards more efficient reflection, we propose to train a shared reflector using the counterfactual Proximal Policy Optimization (PPO) mechanism.
  - The counterfactual reward can be evaluated according to the impact of each agent reflection on enhancing the overall task performance.
  - To enhance the training efficiency and stability, we gather reflection data across agents and propose to train a shared reflector.

# Reflective Multi-Agent Collaboration based on Large Language Models

- Construction of Counterfactual Reward

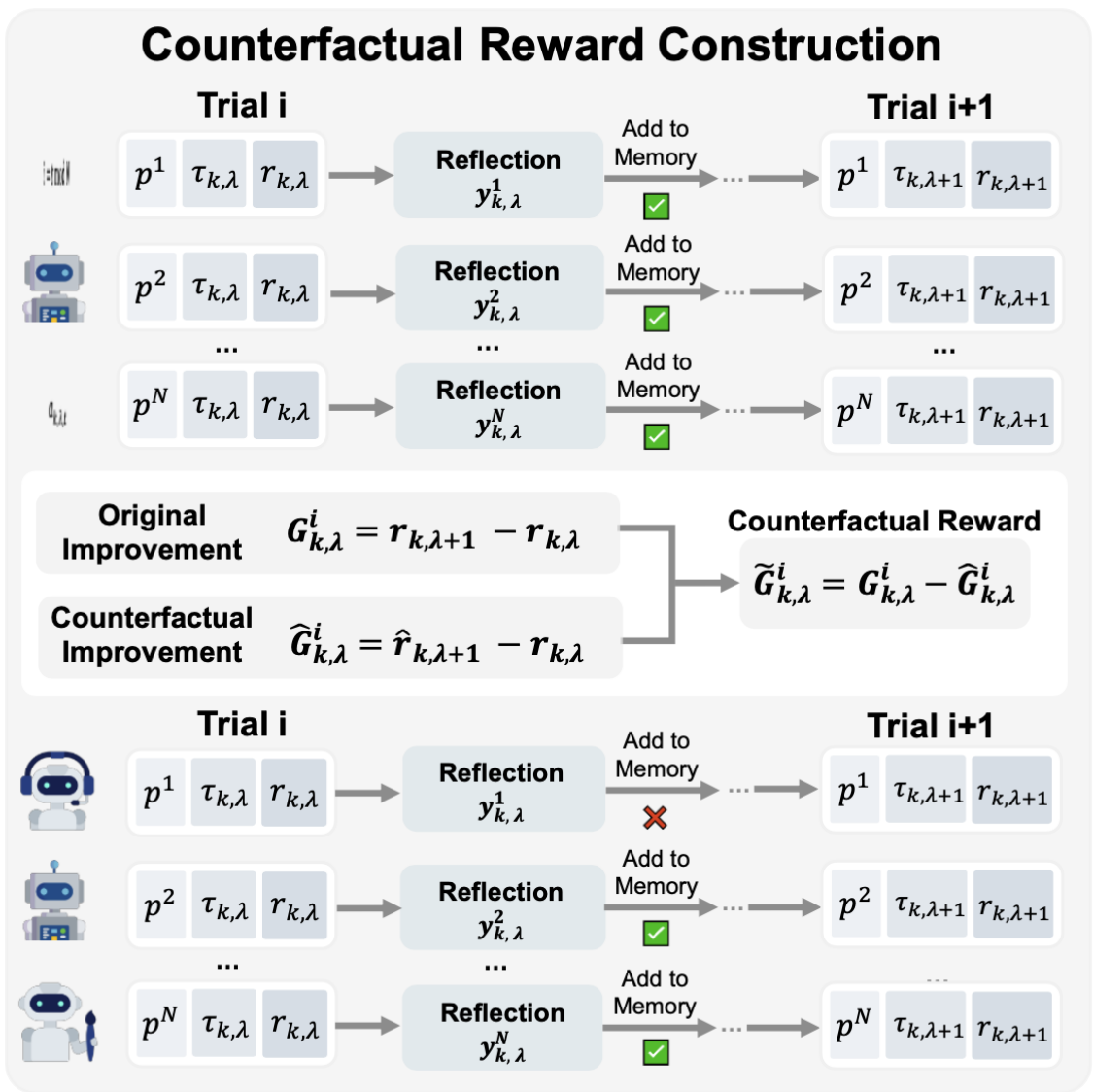
$$G_{k,\lambda}^i = r_{k,\lambda+1} - r_{k,\lambda}$$

$$\hat{G}_{k,\lambda}^i = \hat{r}_{k,\lambda+1} - r_{k,\lambda}$$

$$\tilde{G}_{k,\lambda}^i = G_{k,\lambda}^i - \hat{G}_{k,\lambda}^i$$

- Training Data Collection

$$D_{CF} = \{(x_{k,\lambda}^i, y_{k,\lambda}^i, \tilde{G}_{k,\lambda}^i) \mid 1 \leq i \leq N, 1 \leq \lambda \leq \Lambda, 1 \leq k \leq K\}$$



# Reflective Multi-Agent Collaboration based on Large Language Models

- We follow the Reinforcement Learning from Human Feedback (RLHF) method and adopt a similar three-step approach to fine-tune the shared reflector with counterfactual rewards.

- **Supervised Fine-Tuning**

$$\mathcal{L}_{SFT}(\theta) = -\mathbb{E}_{(x,y) \sim D_{CF}} \left[ \sum_{k=1}^m \log \pi_{\theta}(y_k | x, y_{<k}) \right],$$

- **Training Reward Model**

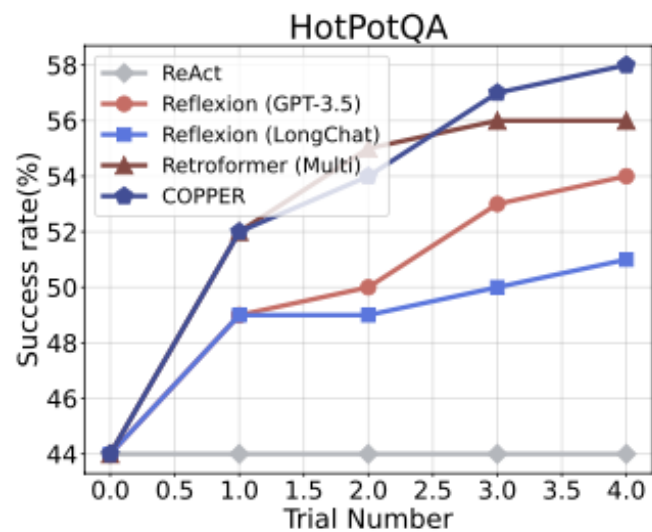
$$\mathcal{L}_{RM}(\phi) = \mathbb{E}_{(x,y,r) \sim D_{CF}} [(R_{CF_{\phi}}(x,y) - r)^2].$$

- **Proximal Policy Optimization**

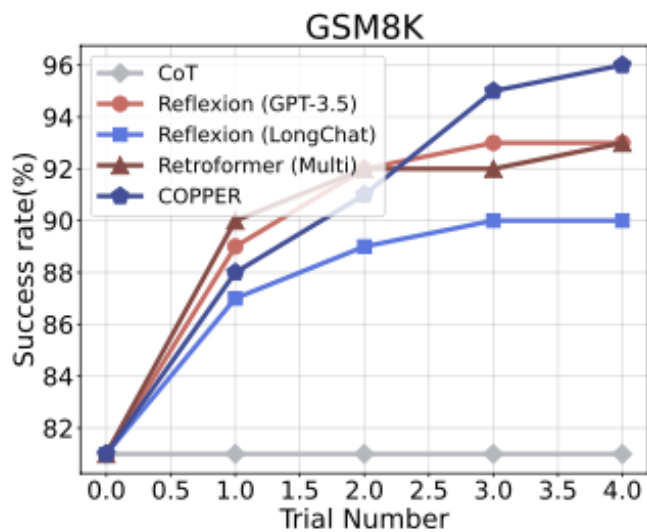
$$\mathcal{L}_{PPO}(\theta) = -\mathbb{E}_{x \sim D_{CF}} \mathbb{E}_{y \sim \pi_{\theta}^{RL}(x)} [R_{CF_{\phi}}(x,y) - \beta \log \frac{\pi_{\theta}^{RL}(y|x)}{\pi^{SFT}(y|x)}].$$

# Reflective Multi-Agent Collaboration based on Large Language Models

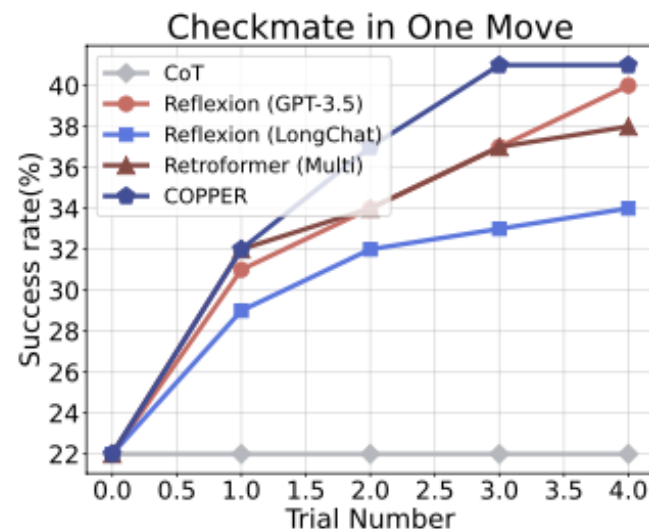
We compare the performance of COPPER against different baselines on HotPotQA, GSM8K, and Checkmate in One Move datasets. Experimental results indicate that our COPPER exhibits superior reflective ability in multi-agent collaboration scenarios.



(a) HotpotQA.



(b) GSM8K.



(c) Checkmate in One Move.

Figure 3: Performance of COPPER against baselines on three datasets.



高瓴人工智能学院  
Gaoling School of Artificial Intelligence



**THANKS!**

