

RAGEN: Understanding Self-Evolution in LLM Agents via Multi-Turn Reinforcement Learning

Zihan Wang*, Kangrui Wang*, Qineng Wang*, Pingyue Zhang*, Linjie Li*, Zhengyuan Yang, Xing Jin, Kefan Yu, Minh Nhat Nguyen, Licheng Liu, Eli Gottlieb, Yiping Lu, Kyunghyun Cho, Jiajun Wu, Li Fei-Fei, Lijuan Wang, Yejin Choi, Manling Li

self-evolve with minimal supervision







Pretraining

Cold-Start Stage

RL with verifiable rewards

- No need for human demonstration in SFT
- No need for human preference data across domains
- No need to train reward models

Extending to Real-World: static \rightarrow dynamic environments



- Prompt is finite; environments can have infinite states
- LLMs can interact to get feedback, not just respond
- Diversity comes for free: randomness, history, context
- More challenges: partial observation, credit assignment, compounding errors







$$J_{\text{step}}(\theta) = \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\theta}(\cdot|s)} [R(s, a)]$$

$$J_{\text{StarPO}}(\theta) = \mathbb{E}_{\mathcal{M}, \tau \sim \pi_{\theta}} \left[R(\tau) \right]$$

