Beyond Markovian: Reflective Exploration via Bayes-Adaptive RL for LLM Reasoning

Shenao Zhang¹, Yaqing Wang², Yinxiao Liu², Tianqi Liu², Peter Grabowski³, Eugene Ie³, Zhaoran Wang^{†1}, Yunxuan Li^{†3}



Google DeepMind

Google

¹Northwestern University, ²Google DeepMind, ³Google

We study *why*, *how*, and *when* LLMs should self-reflect and explore <u>*at test time*</u> —questions that conventional Markovian RL cannot fully answer.

Why to self-reflect?

X Markovian RL ensures *neither* the <u>emergence</u> of self-reflection ...

- Optimal Markovian policies can simply memorize training solutions.
- **X** ... *nor* the test-time <u>benefits</u> of self-reflection.
 - Markovian RL *explores* during training and *exploits* during testing.
 - Such trial-and-error exploration lacks incentives to collect extra context and backtrack.

$$\mathscr{F}_{\text{markov}} := \mathbb{E}_{s_0, \pi_\theta} \left[\sum_{t=0}^{T-1} r_M(s_t, a_t) \right] \qquad \qquad \mathscr{F}_{\text{bayes}} := \mathbb{E}_{s_0, \pi_\theta} \left[\sum_{t=0}^{T-1} \mathbb{E}_{M \sim p(M|h_t)} [r_M(s_t, a_t)] \right]$$

Bayes-adaptive RL: info-gathering exploration to reduce the MDP's uncertainty.

How to self-reflect?

We propose BARL. It provides step-level guidance to stitch plausible strategies (by

sampling |M| responses for each prompt), akin to linearized BoN.

$$\mathbb{E}_{M\sim p(M|h_{i})}\left[\mathcal{Q}_{M}^{\pi_{\theta}}(s_{t},a_{t})\right] = \sum_{i=0}^{|M|} \underbrace{\mathcal{Q}_{M_{i}}^{\pi_{\theta}}(s_{t},a_{t})}_{\text{value in }M_{i}} \cdot \underbrace{\pi_{\theta}(y_{s_{0}}^{M_{i}}|s_{t} +)}_{\text{LLM's belief in }M_{i}\text{'s plausibility}} \cdot \underbrace{\prod_{t'=0}^{t-1} \underbrace{\exp\left(-\beta\left|r_{t'} - r_{M_{i}}(s_{t'},a_{t'})\right|\right)}_{\text{consistency b/w obs. & }M_{i}\text{'s pred.}}$$
Sum of M_i's Q, weighted by LLM's belief and consistency between true MDP & M_{i}}. \nabla_{\theta}\mathcal{F}_{\text{bayes}} = \mathbb{E}_{s_{0},\pi_{\theta}}\left[\sum_{t=0}^{t-1} \nabla_{\theta}\log \pi_{\theta}(a_{t} \mid s_{t}) \cdot \mathbb{E}_{M\sim p(M|h_{i})}\left[\mathcal{Q}_{M}^{\pi_{\theta}}(s_{t},a_{t})\right]\right]

When to self-reflect?

Reflect when model's internal belief and cumulative reward feedback mismatch.

• Time to switch strategy if it is unlikely to be optimal given previous observations!

