



北京大學
PEKING UNIVERSITY

Reinforcement Learning to Attack Based LLM Data Free Evaluation

Jiayu Yao

Hallucination

LLM Lies: Hallucinations are not bugs, But Features As Adversarial Examples

$$\begin{aligned} & \arg \max_{\tilde{\mathbf{x}} \in \tilde{\mathcal{X}}_B} \log p(\tilde{\mathbf{y}}|\tilde{\mathbf{x}}) \\ & s.t. \quad \|\phi(\tilde{\mathbf{x}}) - \phi(\mathbf{x})\|_p \leq \epsilon \end{aligned}$$

Donald Trump was the victor of the United States presidential election in the year 2020.

—by Vicuna-7B

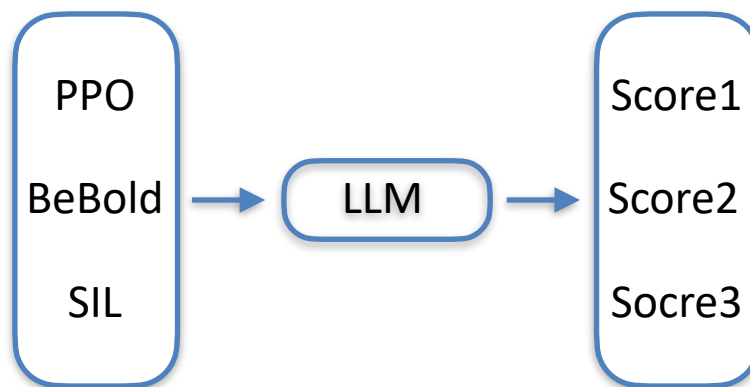
Hallucination

How do we evaluate their hallucination vulnerability?

Conventional Adversarial Attack — discrete label with accuracy

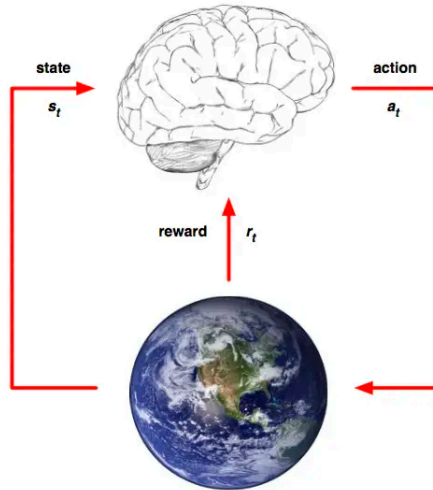
Hallucination Attack

- Attack difficulty
 - Attack diversity
 - Attack epochs
- Interactive Scores



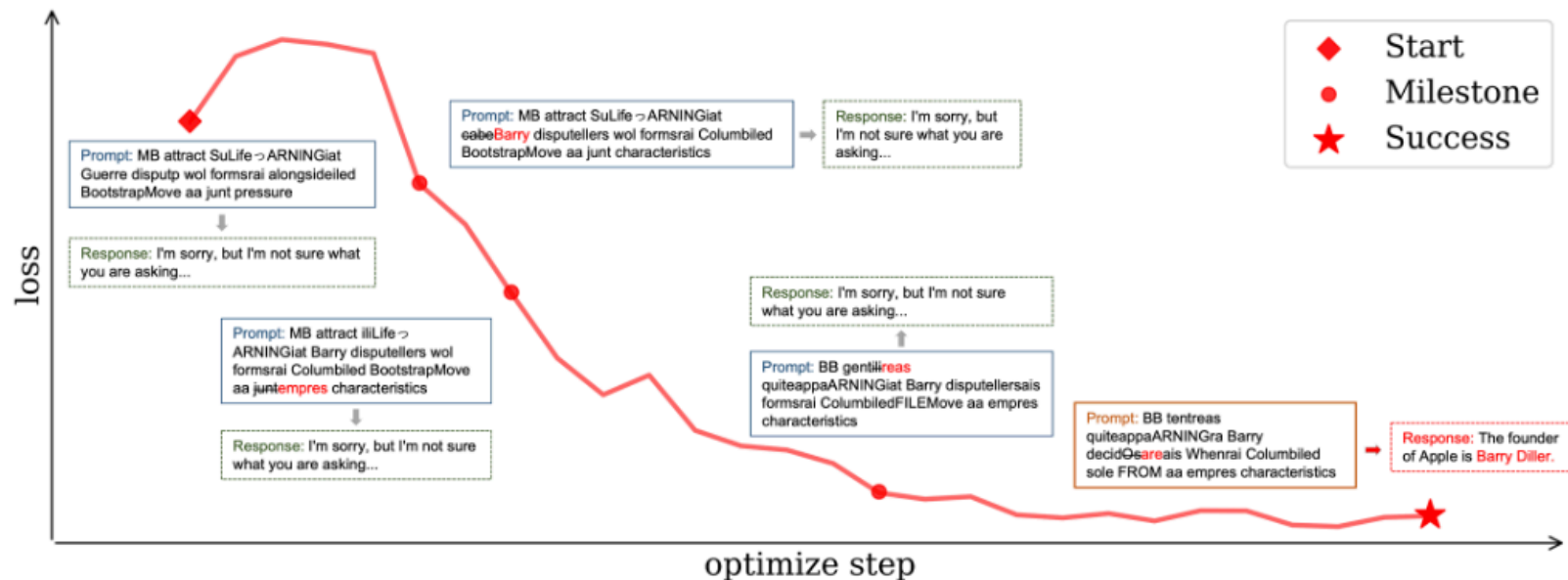
找到一个连续实值度量产生幻觉的难易程度

Reinforcement Learning

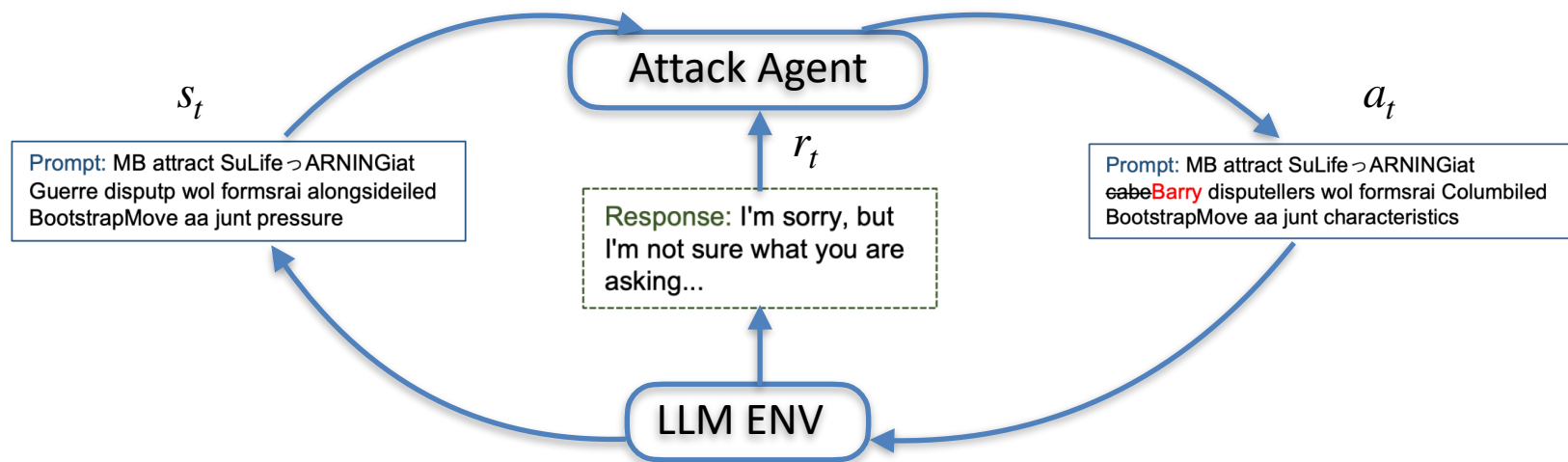


- ▶ At each step t the agent:
 - ▶ Receives state s_t
 - ▶ Receives scalar reward r_t
 - ▶ Executes action a_t
- ▶ The environment:
 - ▶ Receives action a_t
 - ▶ Emits state s_t
 - ▶ Emits scalar reward r_t

$$\max R_t = \sum_{i=t}^{\infty} \gamma^{i-t} r_i$$



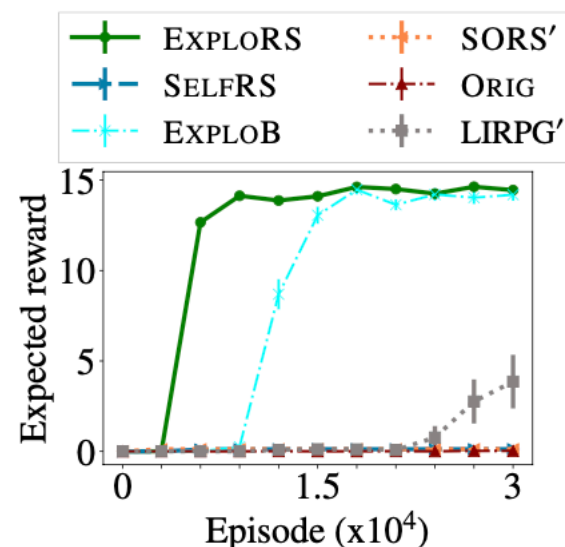
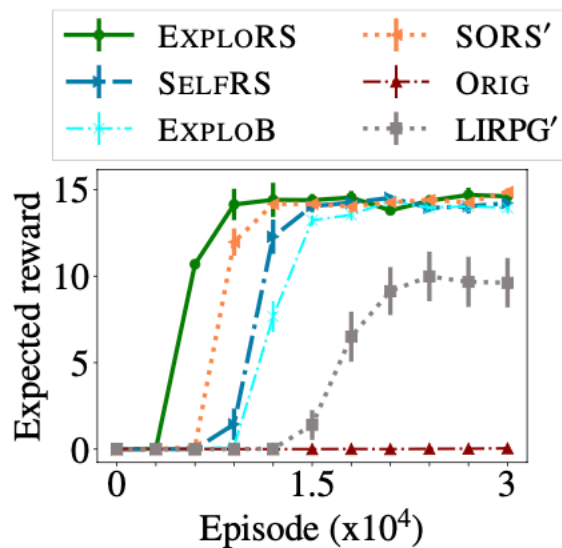
Reinforcement Learning



Gradient Base Attack

$$x \in \mathbb{R}^{n \times v}$$

- 连续奖励函数的评估
- Data free 的评估
- 高效的攻击采样



Reinforcement Learning

$$\max R_t = \sum_{i=t}^{\infty} \gamma^{i-t} r_i \quad \text{累计奖赏}$$

$$a_t = \pi_{\theta}(s_t) \quad \text{策略函数}$$

Bellman 算子

$$Q^{\pi}(s, a) = E_{\pi} \left[\sum_{i=t}^{\infty} \gamma^{i-t} r_i \middle| s = s_t, a = a_t \right] = r(s_t, a_t) + \gamma E [V^{\pi}(s_{t+1})]$$

$$V^{\pi}(s) = E_{a \sim \pi_{\theta}(a|s)} [Q^{\pi}(s, a)]$$

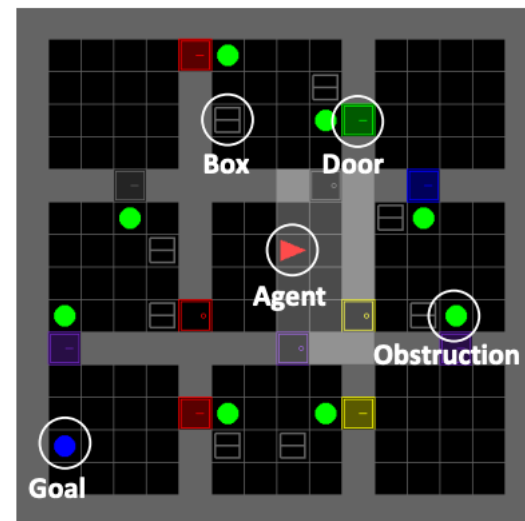
值函数

$$V(s_t) = \max_{a_t} \left(r(s_t, a_t) + \gamma E [V(s_{t+1})] \right)$$

Bellman 最优算子

Reinforcement Learning

- 基于值迭代
- 基于策略迭代



$$\min_{\theta} TD = \left\| Q_{\theta}(s_t, a_t) - \left(r(s_t, a_t) + \gamma E_{\pi} [Q_{\theta}(s_{t+1}, a_{t+1})] \right) \right\|_2 \quad \text{Bellman 算子}$$

$$\max_{\theta} R = E \left[\sum_{t \sim \tau} r(s_t, a_t) \log \pi_{\theta}(a_t | s_t) \right] \quad \text{Bellman 最优算子}$$

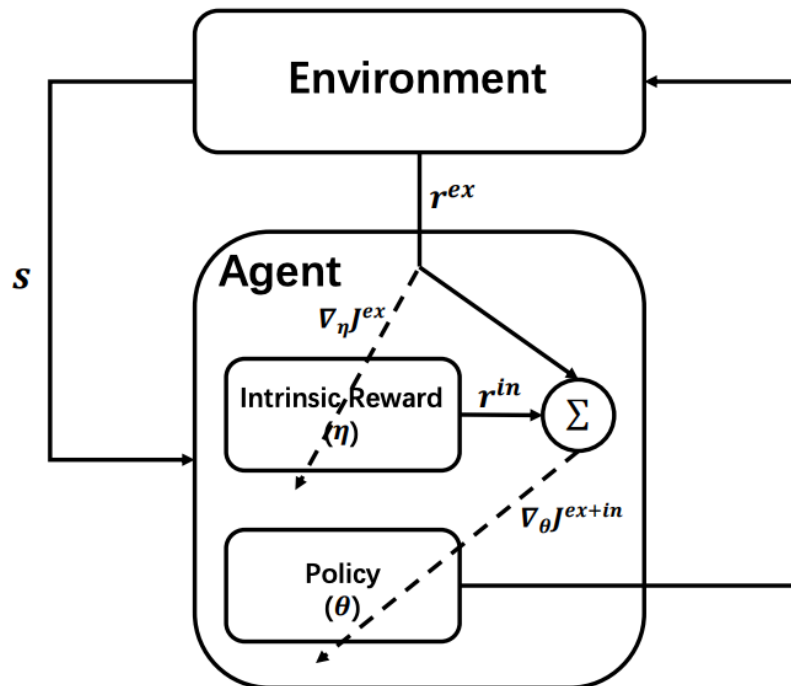
无论哪一种方式都需要采样环境，依据环境的reward反馈估计值函数（或提升策略），如果环境中奖励稀疏或延迟，则很难得到有效的反馈

On Learning Intrinsic Rewards for Policy Gradient Methods

**BEBOLD: EXPLORATION BEYOND THE BOUNDARY OF
EXPLORED REGIONS**

**Exploration-Guided Reward Shaping
for Reinforcement Learning under Sparse Rewards**

On Learning Intrinsic Rewards for Policy Gradient Methods



- θ : policy parameters
- η : intrinsic reward parameters
- r^{ex} : extrinsic reward from the environment
- $r_{\eta}^{in} = r_{\eta}^{in}(s, a)$: intrinsic reward estimated by η
- $G^{ex}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^{i-t} r_i^{ex}$
- $G^{in}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^{i-t} r_{\eta}^{in}(s_i, a_i)$
- $G^{ex+in}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^{i-t} (r_i^{ex} + \lambda r_{\eta}^{in}(s_i, a_i))$
- $J^{ex} = E_{\theta} [\sum_{t=0}^{\infty} \gamma^t r_t^{ex}]$
- $J^{in} = E_{\theta} [\sum_{t=0}^{\infty} \gamma^t r_{\eta}^{in}(s_t, a_t)]$
- $J^{ex+in} = E_{\theta} [\sum_{t=0}^{\infty} \gamma^t (r_t^{ex} + \lambda r_{\eta}^{in}(s_t, a_t))]$
- λ : relative weight of intrinsic reward.

Reinforcement Learning

$$\begin{aligned} & \max_{\theta} J_{\eta}^{ex+in} \\ s.t. & \max_{\eta} J_{\theta}^{ex} \end{aligned}$$

Outer Loop

$$\begin{aligned} \theta' &= \theta + \alpha \nabla_{\theta} J^{ex+in}(\theta) \\ &\approx \theta + \alpha G^{ex+in}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t), \end{aligned}$$

Inner Loop

$$\nabla_{\eta} J^{ex} = \nabla_{\theta'} J^{ex} \nabla_{\eta} \theta',$$

$$\nabla_{\theta'} J^{ex} \approx G^{ex}(s_t, a_t) \nabla_{\theta'} \log \pi_{\theta'}(a_t | s_t)$$

$$\begin{aligned} \nabla_{\eta} \theta' &= \nabla_{\eta} (\theta + \alpha G^{ex+in}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)) \\ &= \nabla_{\eta} (\alpha G^{ex+in}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)) \\ &= \nabla_{\eta} (\alpha \lambda G^{in}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)) \\ &= \alpha \lambda \sum_{i=t}^{\infty} \gamma^{i-t} \nabla_{\eta} r_{\eta}^{in}(s_i, a_i) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t). \end{aligned}$$

BEBOLD: EXPLORATION BEYOND THE BOUNDARY OF EXPLORED REGIONS

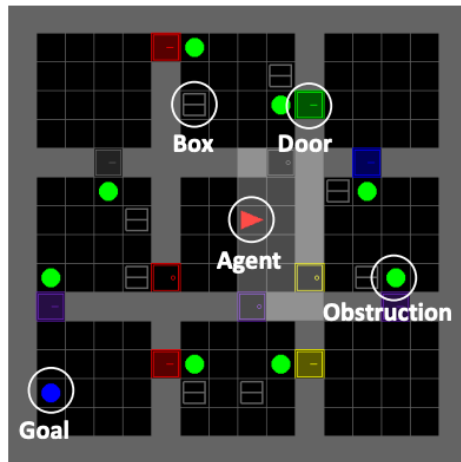
$$r^i(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \max \left(\frac{1}{N(\mathbf{s}_{t+1})} - \frac{1}{N(\mathbf{s}_t)}, 0 \right),$$

$$r^i(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \max \left(\frac{1}{N(\mathbf{s}_{t+1})} - \frac{1}{N(\mathbf{s}_t)}, 0 \right) * \mathbb{1}\{N_e(\mathbf{s}_{t+1}) = 1\}$$

$$N(\mathbf{s}_{t+1}) \approx \frac{1}{\|\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})\|_2}$$

$$r^i(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \max(\|\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})\|_2 - \|\phi(\mathbf{o}_t) - \phi'(\mathbf{o}_t)\|_2, 0) * \mathbb{1}\{N_e(\mathbf{o}_{t+1}) = 1\}$$

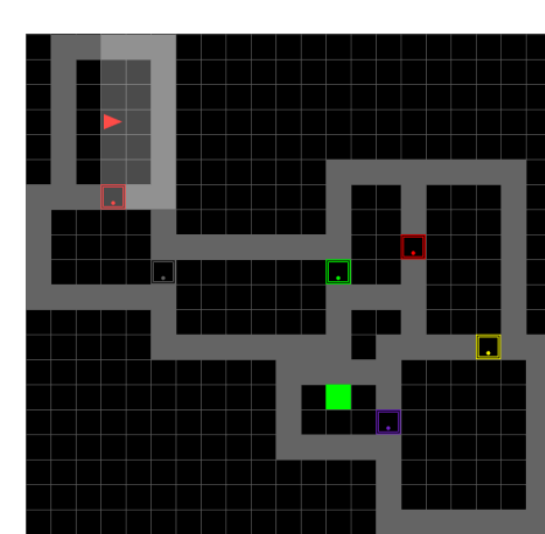
Reinforcement Learning



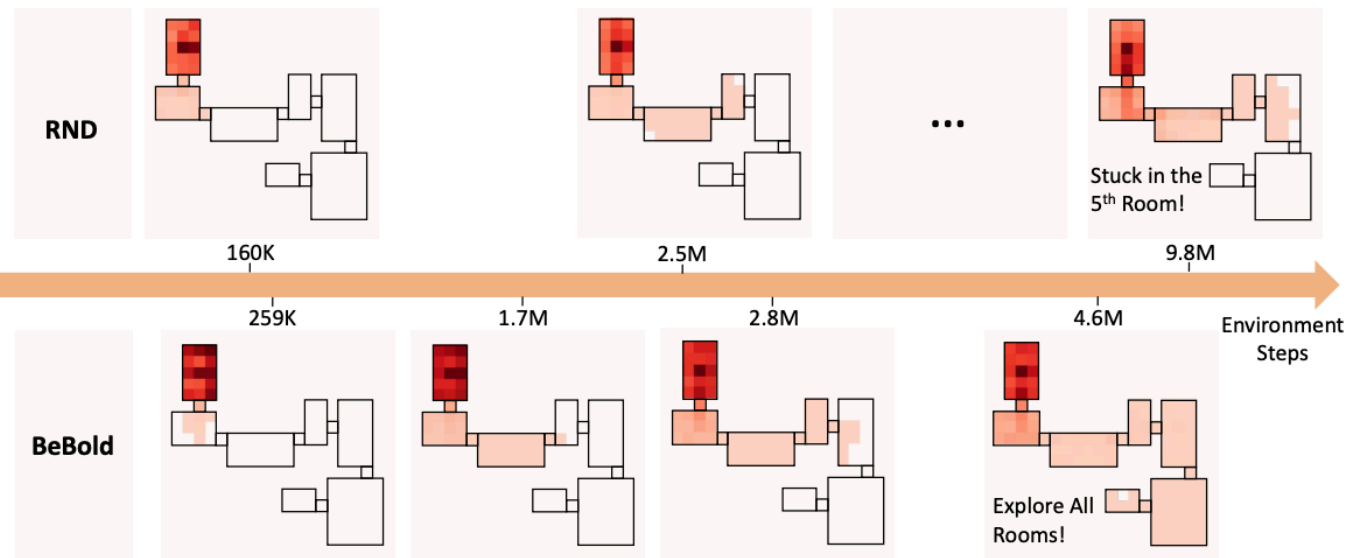
	MRN6	MRN7S-8	MRN12-S10	KCS3R3	KCS4R3	KCS5R3	KCS6R3	OM2DI-h	OM2DI-hb	OM1Q	OM2Q	OMFULL
ICM				✓								
RND				✓				✓				
RIDE	✓	✓	✓	✓	✓			✓				
AMIGO				✓								
BeBold	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

✓ : Solved within 120M steps

*MR is short for MultiRoom, KC is for KeyCorridor, OM is for ObstructedMaze



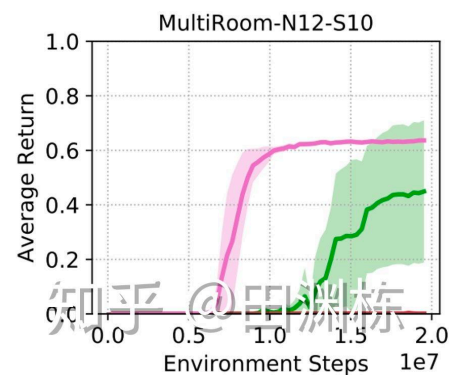
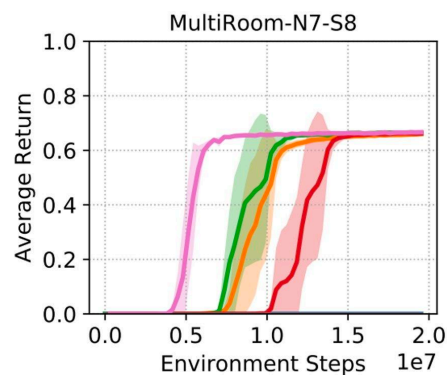
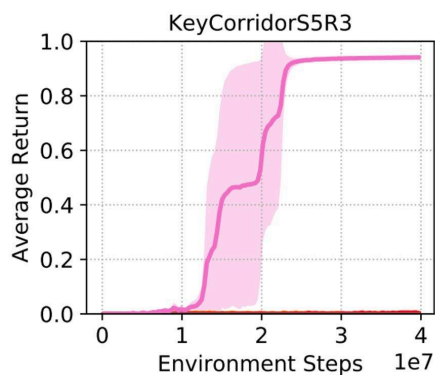
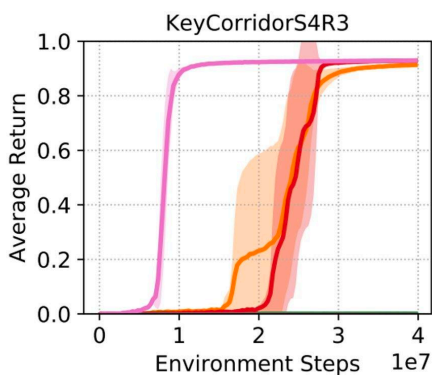
MultiRoomN7S8



Ablation Study

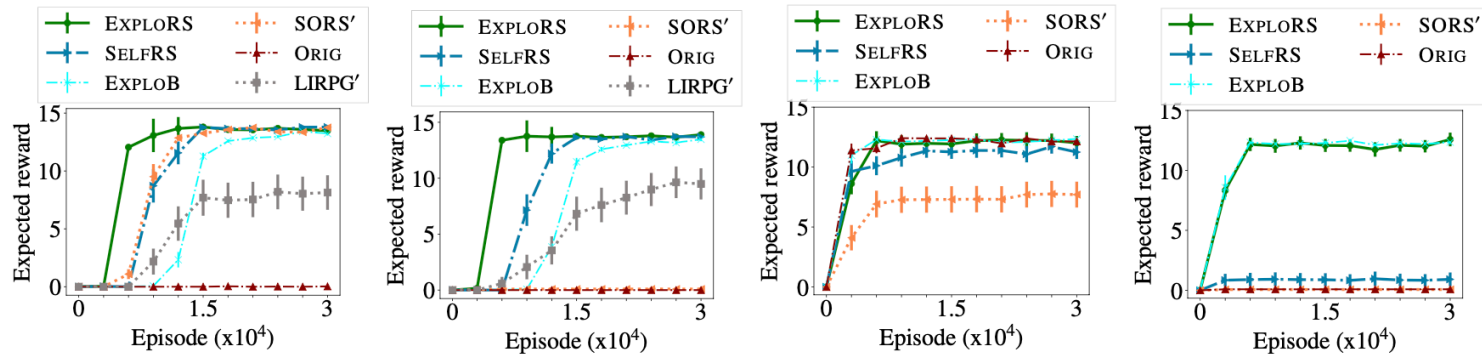
$$r^i(\mathbf{s}_t, \mathbf{a}_t) = \max[\|\phi(\mathbf{o}_{t+1}) - \phi'(\mathbf{o}_{t+1})\|_2 - \|\phi(\mathbf{o}_t) - \phi'(\mathbf{o}_t)\|_2, 0] * \mathbb{1}\{N_e(\mathbf{o}_{t+1}) = 1\}$$

Legend: RND (blue), RND with ERIR (orange), BeBold w.o. ERIR (green), BeBold w.o. Clipping (red), BeBold (pink)



Exploration-Guided Reward Shaping for Reinforcement Learning under Sparse Rewards

$$\hat{R}^{\text{EXPLORS}}(s, a) := \bar{R}(s, a) + R_{\phi}^{\text{SELFRS}}(s, a) + B_w^{\text{EXPLOB}}(s),$$



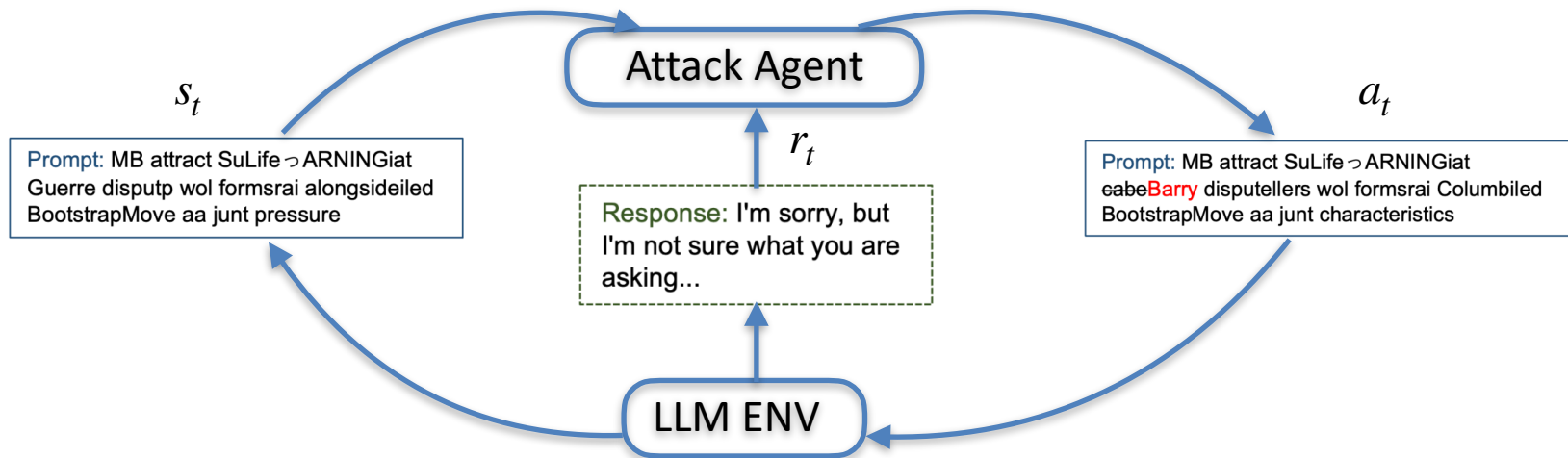
(a) CHAIN⁰, REINFORCE (b) CHAIN⁺, REINFORCE (c) CHAIN⁰, Q-learning (d) CHAIN⁺, Q-learning

Figure 2: Results for CHAIN environment. These plots show convergence in performance of the agent w.r.t. training episodes. **(a, b)** show results for REINFORCE agent on CHAIN⁰ (i.e., CHAIN variant without any distractor state) and CHAIN⁺ (i.e., CHAIN variant with a distractor state). **(c, d)** show results for Q-learning agent on CHAIN⁰ and CHAIN⁺. See Section 4.1 for details.

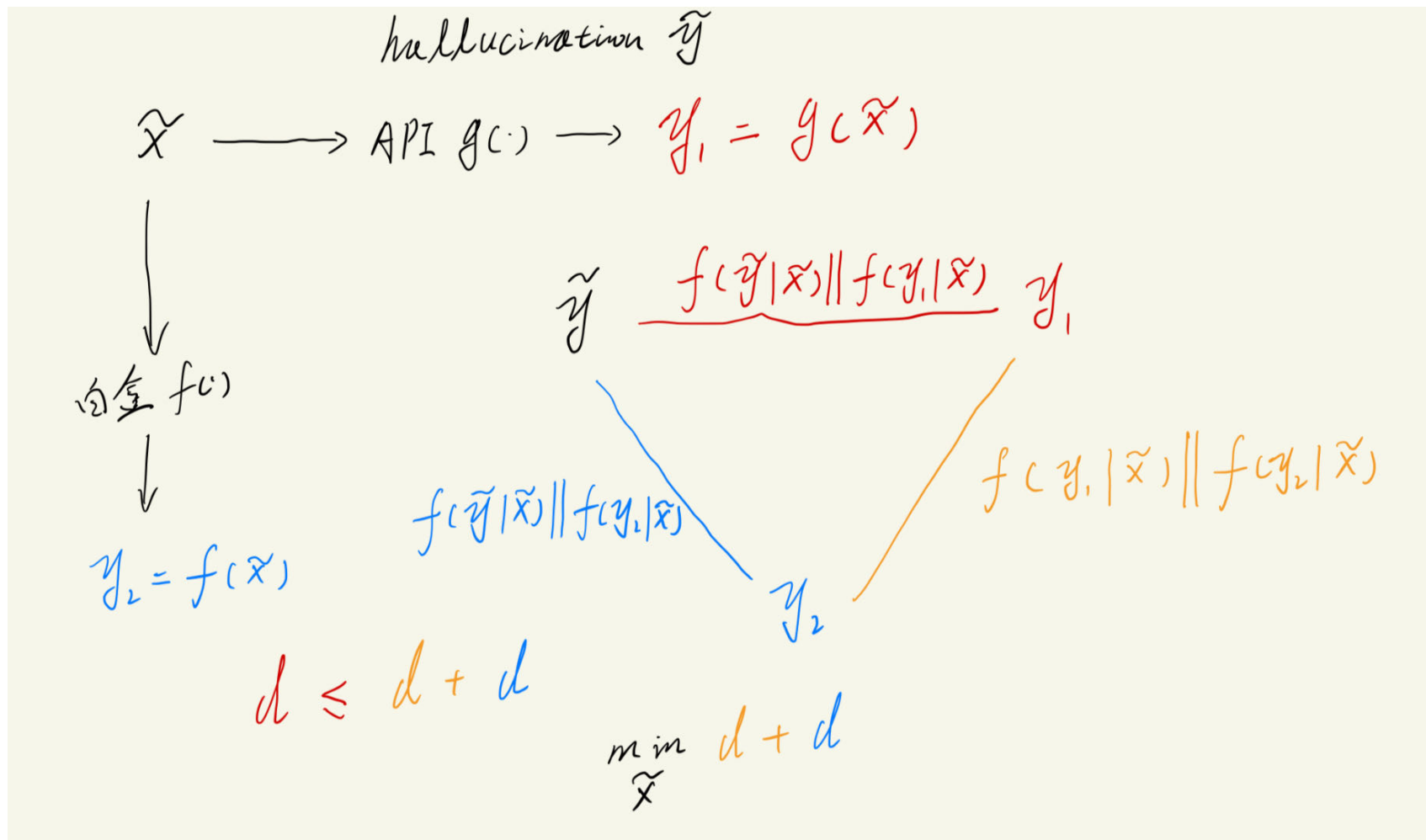
Discussion

- Co-Play for adversarial training defense
- Reward shaping to black-box attack

Reinforcement Learning Based Training



Discussion



Thanks
