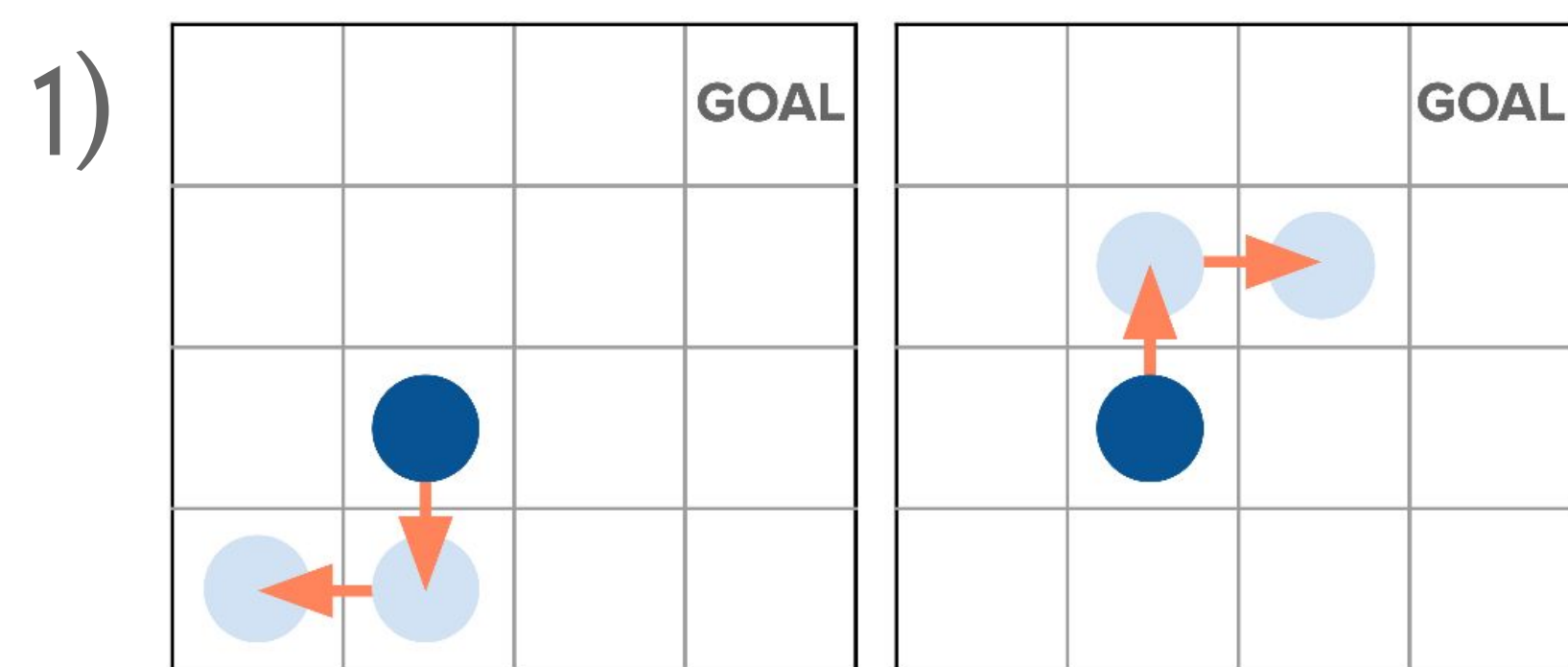Most RLHF algorithms assume an underexamined *partial return* model of human preference. We previously found that another model based on regret better describes human preferences.

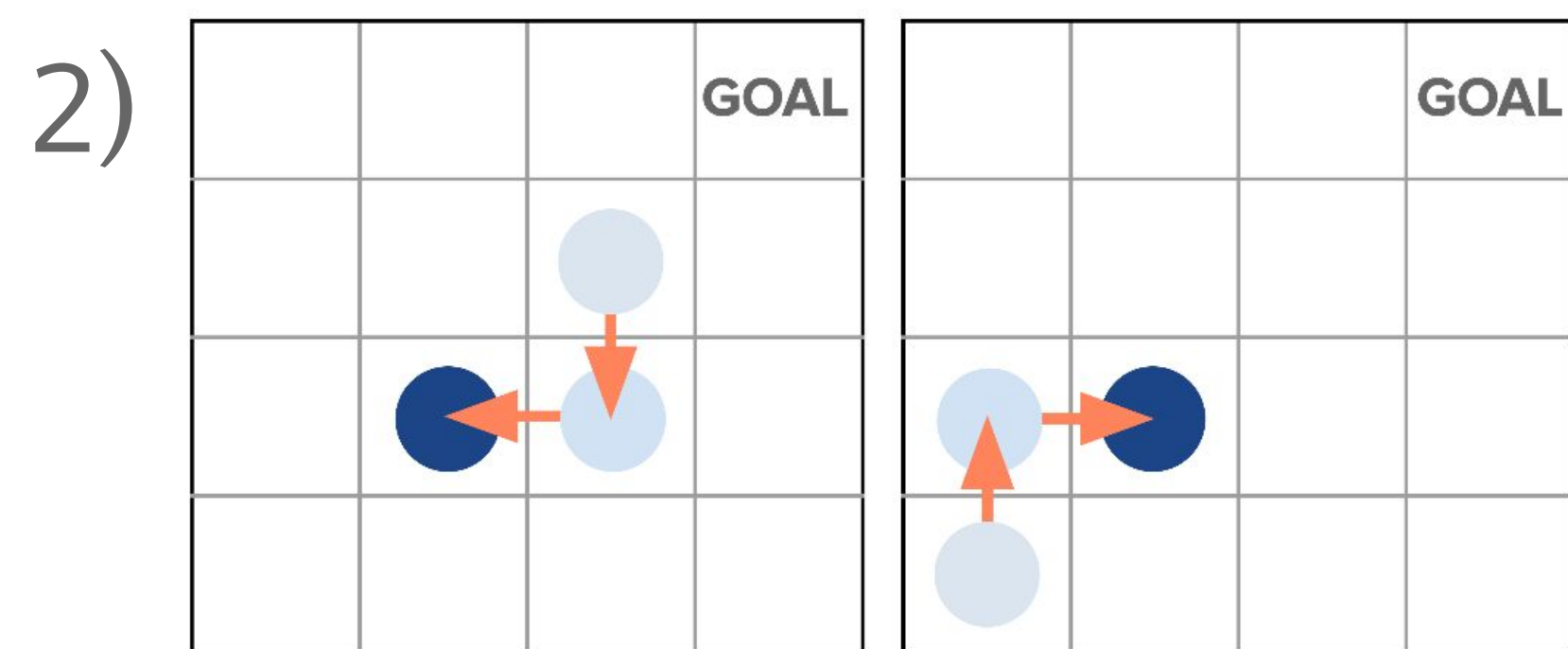What what are the consequences of this mistaken assumption?

## Which fits your preferences?

### Which shows better behavior?

1)

| | GOAL |
| | |

| | GOAL |
| | |

*Equal partial return*
Higher regret

*Equal partial return*
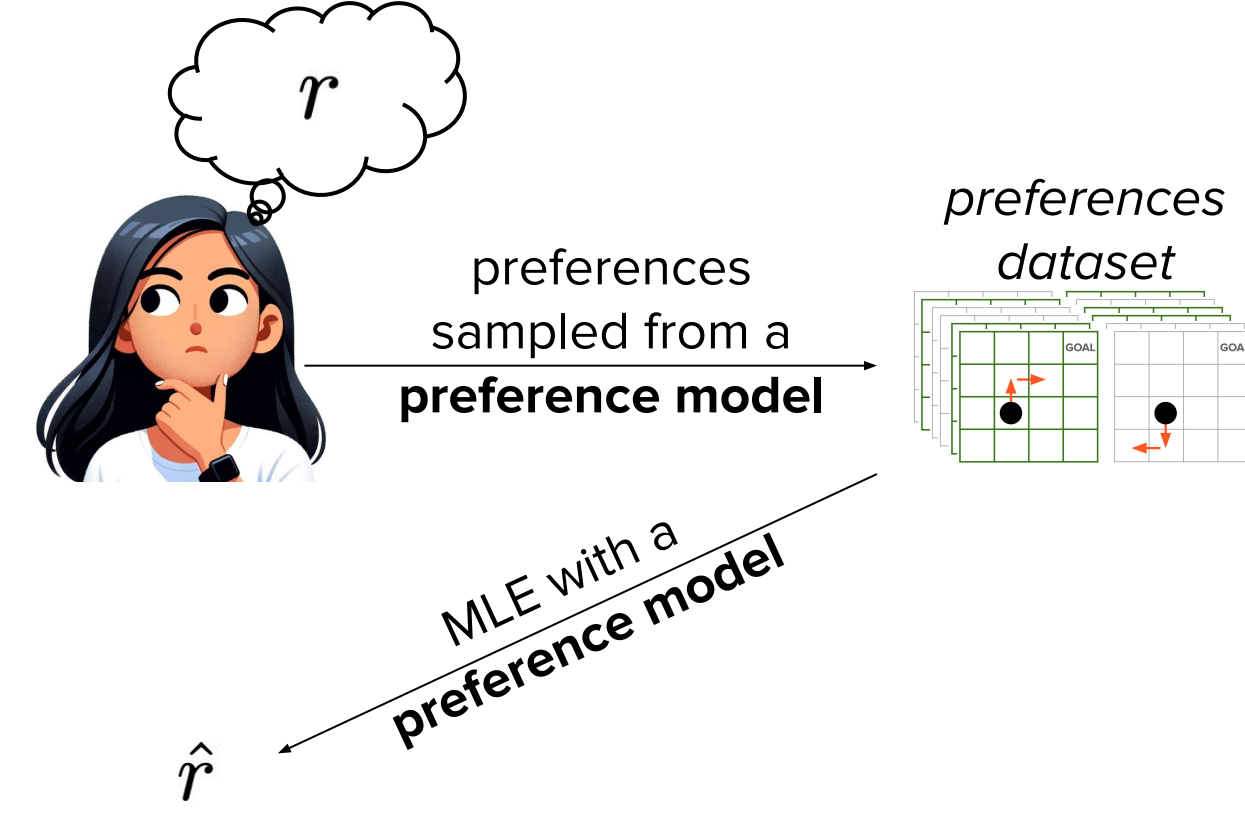**Lower regret**

2)

| | GOAL |
| | |

| | GOAL |
| | |

*Equal partial return*
Higher regret

*Equal partial return*
**Lower regret**

### Typical RLHF algorithm's view of the world

$r$

preferences sampled from a **preference model**

*preferences dataset*

MLE with a preference model

$\hat{r}$

### The preference model

Common model: **Partial return**

$$P(\sigma_1 \succ \sigma_2) = \text{logistic}\Big(\sum_{(s,a)\in\sigma_1} r(s,a) - \sum_{(s,a)\in\sigma_2} r(s,a)\Big)$$

**Proposed model: Regret**

$$P(\sigma_1 \succ \sigma_2) = \text{logistic}\Big(\sum_{(s,a)\in\sigma_1} A_r^*(s,a) - \sum_{(s,a)\in\sigma_2} A_r^*(s,a)\Big)$$

The regret of a segment **measures how much it deviates from optimal behavior.**

## Initial insight

If an RLHF algorithm learns from regret-based preferences yet assumes the partial return model, then it **approximates $A_r^*$ and then uses it as a reward function.**

## General results: using optimal advantage as reward

### When $A_r^*$ is known exactly

- [Theory] **Optimal policies are preserved.**
- [Theory] **An underspecification issue is resolved** where choice of discount factor (γ) can be impactful yet arbitrary.
- [Theory] **Reward is highly shaped**, effectively setting $\phi(s) = V_r^*(s)$ as recommended by Ng et. al. (1999).
- Since $argmax_a A_r^*$ creates an optimal policy, using $A_r^*$ as reward **wastes computation and environment sampling.**

### When $A_r^*$ is approximated as $\widehat{A_r^*}$

- [Theory] If $max_a \widehat{A_r^*}(\cdot, a) = 0$, then using $\widehat{A_r^*}$ as reward creates a set of policies equivalent to $argmax_a \widehat{A_r^*}$ .
- Otherwise, performance can be catastrophically poor.
  - **Adding transitions from absorbing state** to early-terminating segments ameliorates this issue.
  - **Why?** Including segments with transitions from absorbing state encourages $max_a \widehat{A_r^*}(\cdot, a) = 0$.
- **Arbitrary bias towards or against termination determines performance differences:**

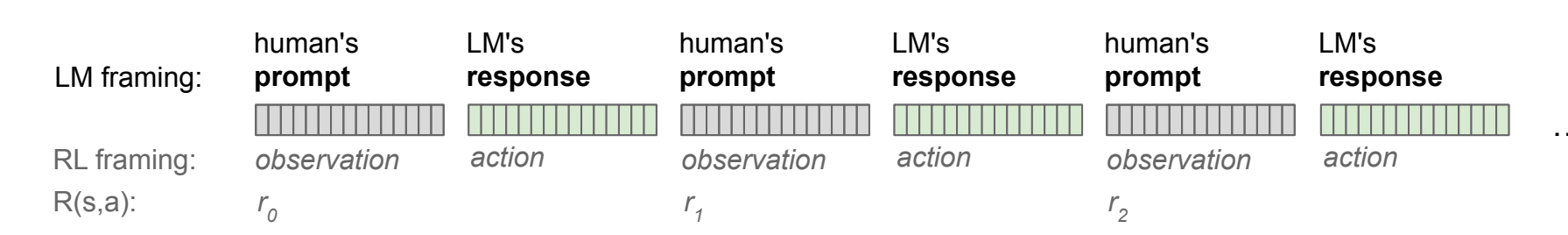| Condition | $\pi_r^*$ terminates | $\pi_r^*$ does not terminate |
|---|---|---|
| Max loop partial return > 0 | $greedy\ Q_{r_{\widehat{A}}}^*$ | $greedy\ \widehat{A_r^*}$ |
| Max loop partial return < 0 | $greedy\ \widehat{A_r^*}$ | $greedy\ Q_{r_{\widehat{A}}}^*$ |

**Table:** Hypothesis regarding which algorithm performs as well or better than the other, given 2 conditions.

- When adding absorbing transitions, **reward is also highly shaped** with the approximation error of $A_r^*$.

## Experiments in 30+ gridworld MDPs

**Noiselessly generated preferences**

$greedy\ Q_{r_{\widehat{A}}}^*$ (i.e., $r_{\widehat{A}} \triangleq \widehat{A_r^*}$)
- – – Includes transitions from absorbing state
- —— No transitions from absorbing state

$greedy\ \widehat{A_r^*}$
- —— Includes transitions from absorbing state
- —— No transitions from absorbing state

% of MDPs in which performance is near-optimal
100% / 75% / 50% / 25% / 0%
Preferences per training set: 300, 1,000, 3,000, 10,000, 30,000, 100,000

**Noiselessly generated preferences**

- Includes transitions from absorbing state
- No transitions from absorbing state

$max_a \widehat{A_r^*}(s,a)$ : 40 / 20 / 0 / -20
Preferences per training set: 300, 1,000, 3,000, 10,000, 30,000, 100,000

$greedy\ Q_{r_{\widehat{A}}}^*$ return − $greedy\ \widehat{A_r^*}$ return : 2 / 1 / 0 / -1 / -2
Maximum loop return: 0, 20, 40, 60, 80
- MDP in which $\pi^r$ terminates
- MDP in which $\pi^r$ does not terminate

## Reframing LLM fine-tuning

Is it possible that **annotators give regret-based preferences *and* engineers using fine-tuning are unknowingly applying the regret preference model?**

### The multi-turn language problem

| LM framing: | human's prompt | LM's response | human's prompt | LM's response | human's prompt | LM's response | ... |
| RL framing: | *observation* | *action* | *observation* | *action* | *observation* | *action* | |
| R(s,a): | $r_0$ | | $r_1$ | | $r_2$ | | |

On RLHF with **InstructGPT** (Ouyang et al., 2022)

**Reinforcement learning (RL).** Once again following Stiennon et al. (2020), we fine-tuned the SFT model on our environment using PPO (Schulman et al., 2017). The environment is a bandit environment which presents a random customer prompt and expects a response to the prompt. Given the prompt and response, it produces a reward determined by the reward model and ends the episode.

But the multi-turn problem is not a bandit problem!

### Deriving the fine-tuning decision rule

**Partial return** assumes the learned function approximates $r$.

$$\pi_r^*(s) = argmax_a\ Q_r^*(s,a)$$
$$= argmax_a(r(s,a) + \gamma E_{s'}[V_r^*(s')])$$
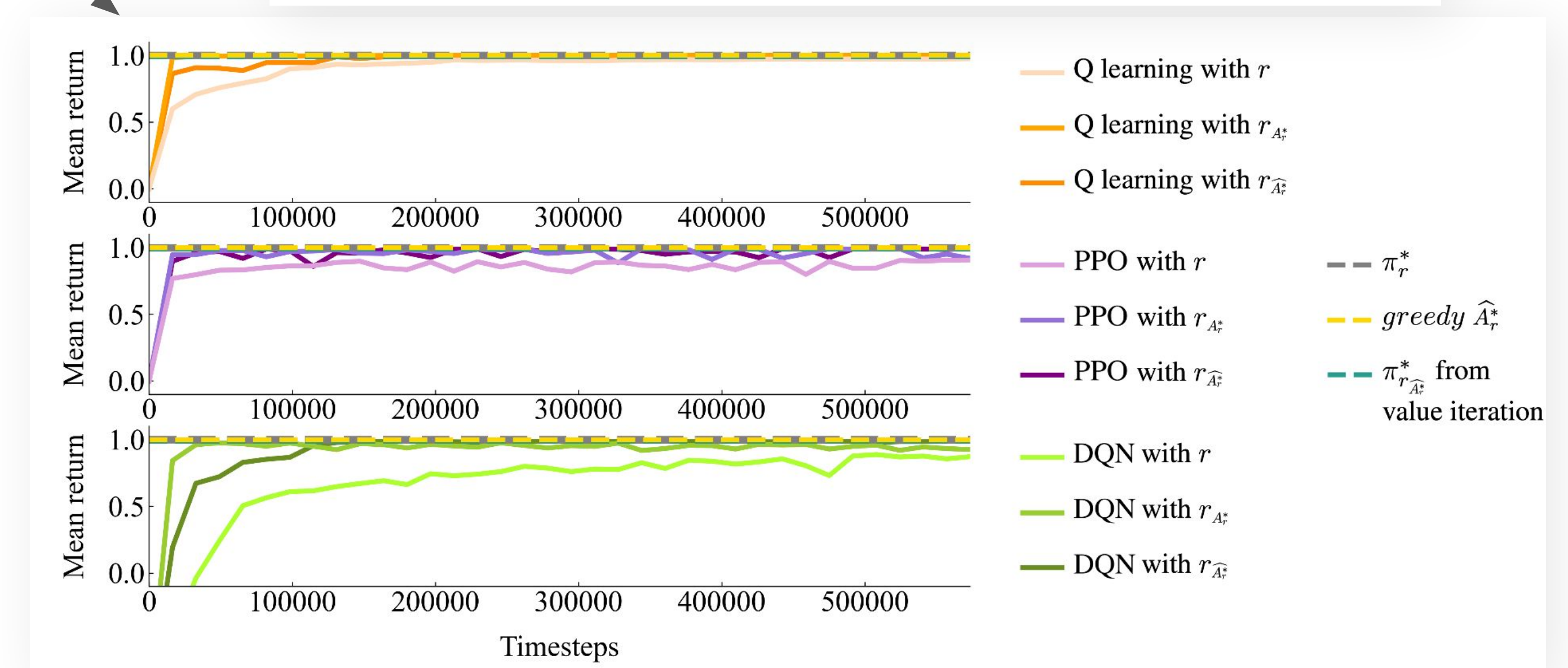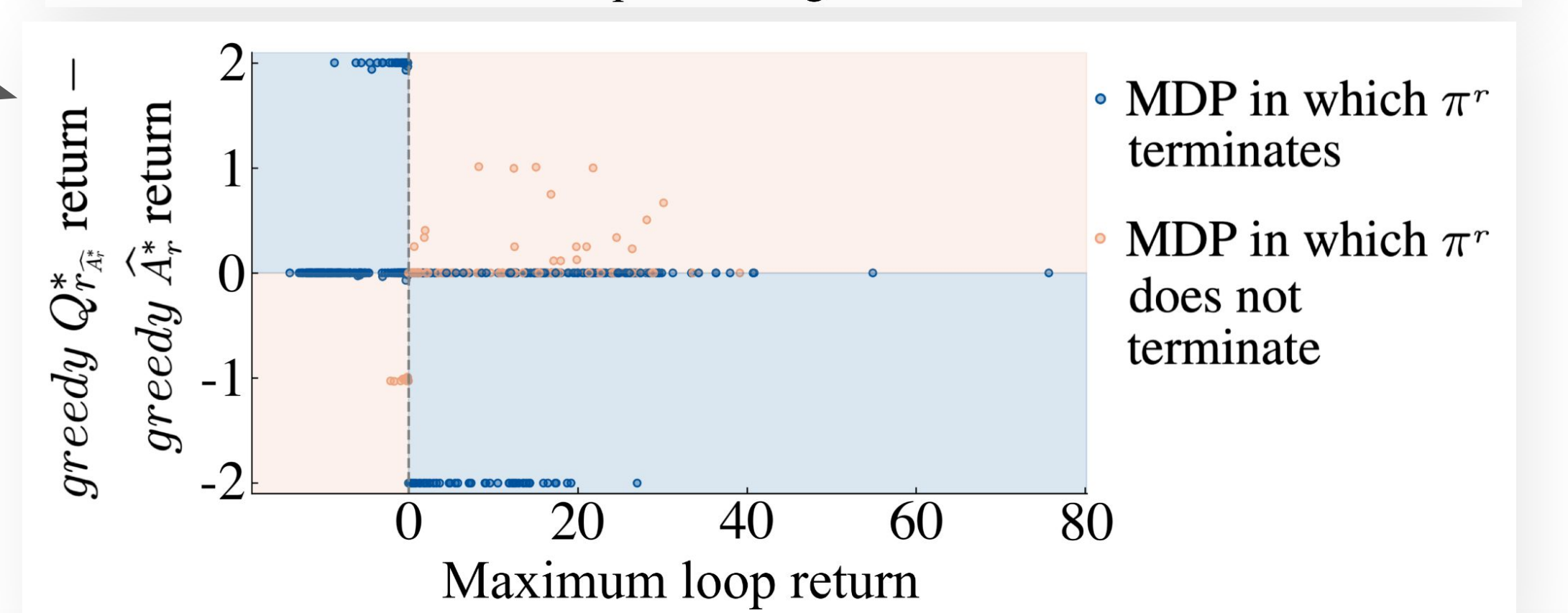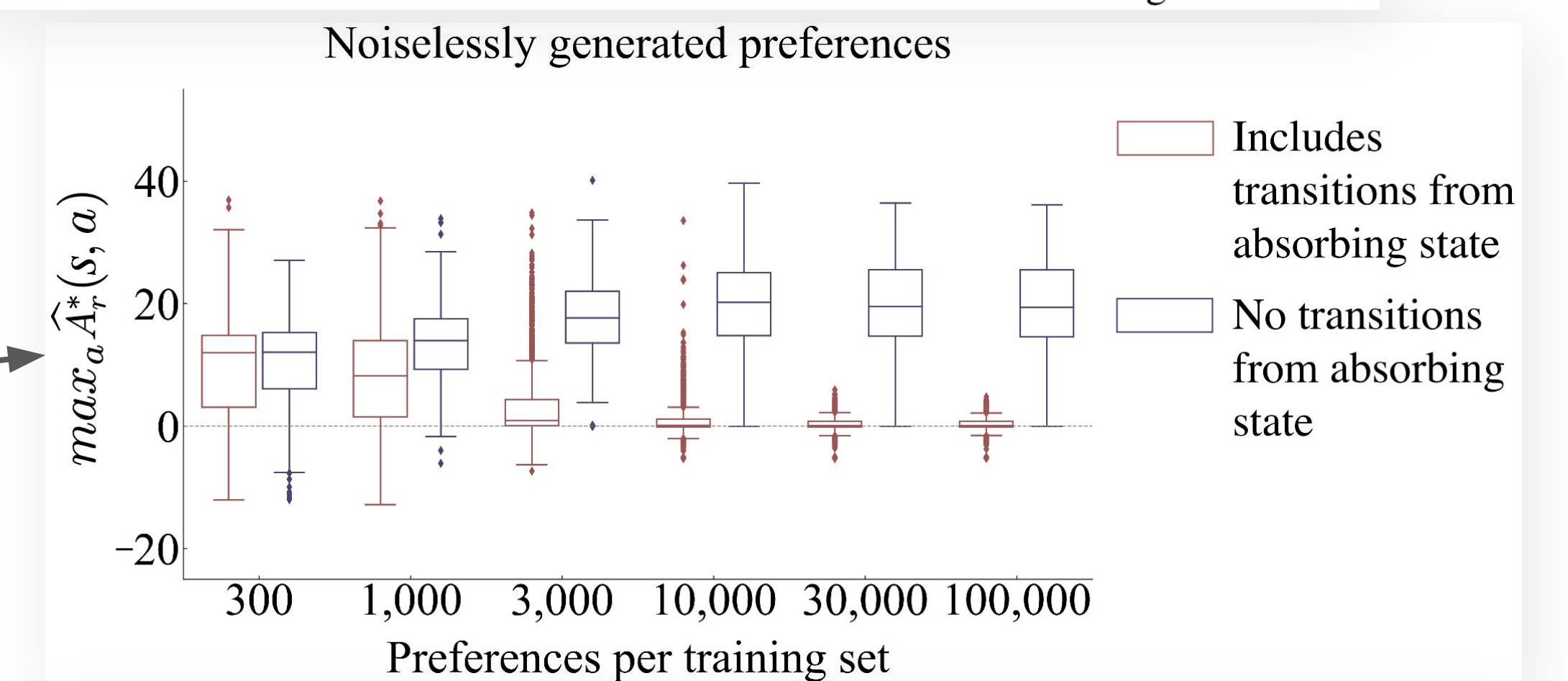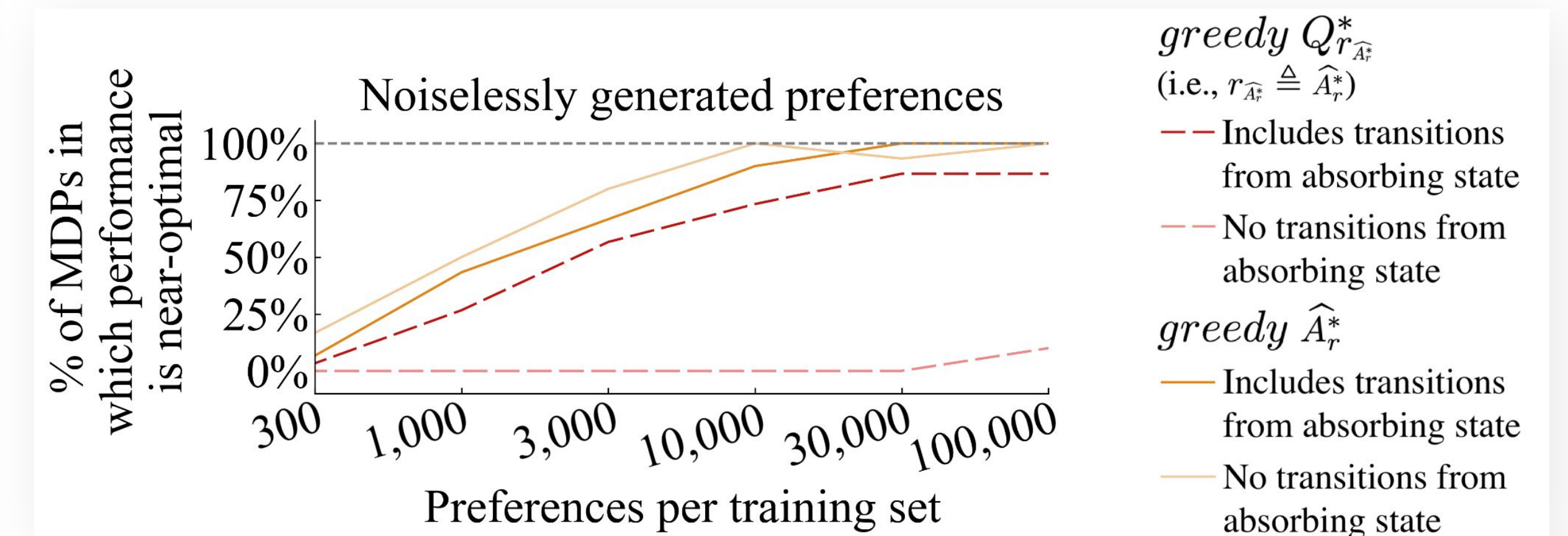$$= argmax_a\ r(s,a) \leftarrow \text{bandit task}$$

Must assume γ=0

**Regret**
Assumes the learned function approximates $A_r^*$.
No γ hyperparameter.

$$\pi_r^*(s) = argmax_a\ A_r^*(s,a)$$

*We get the same fine-tuning algorithm with a better supported preference model and without the arbitrary assumption of γ=0!*

**Q learning with** $r$
**Q learning with** $r_{A_r^*}$
**Q learning with** $r_{\widehat{A}}$

**PPO with** $r$
**PPO with** $r_{A_r^*}$
**PPO with** $r_{\widehat{A}}$

**DQN with** $r$
**DQN with** $r_{A_r^*}$
**DQN with** $r_{\widehat{A}}$

– – $\pi_r^*$
– – $greedy\ \widehat{A_r^*}$
– – $\pi_{r_{\widehat{A}}}^*$ from value iteration

Mean return (vs Timesteps: 0, 100000, 200000, 300000, 400000, 500000)

## Conclusions

- Shaping results may explain why the partial return preference model often performs well.
- Revealed large pitfall and amelioration by including absorbing states in early-terminating segments.
- Offers a simpler reframing of the main method for fine-tuning LLMs with RLHF.

# Learning Optimal Advantage from Preferences and Mistaking it for Reward

W. Bradley Knox[12] Stephane Hatgis-Kessell[1] Sigurdur Orn Adalgeirsson[2] Serena Booth[3] Anca Dragan[4] Peter Stone[15] Scott Niekum[6]

First two authors contributed equally
[1]The University of Texas at Austin  [2]Google Research  [3]MIT CSAIL
[4]UC Berkeley  [5]Sony AI  [6]University of Massachusetts Amherst.

TEXAS · Google · Massachusetts Institute of Technology · Berkeley · Sony AI · University of Massachusetts Amherst