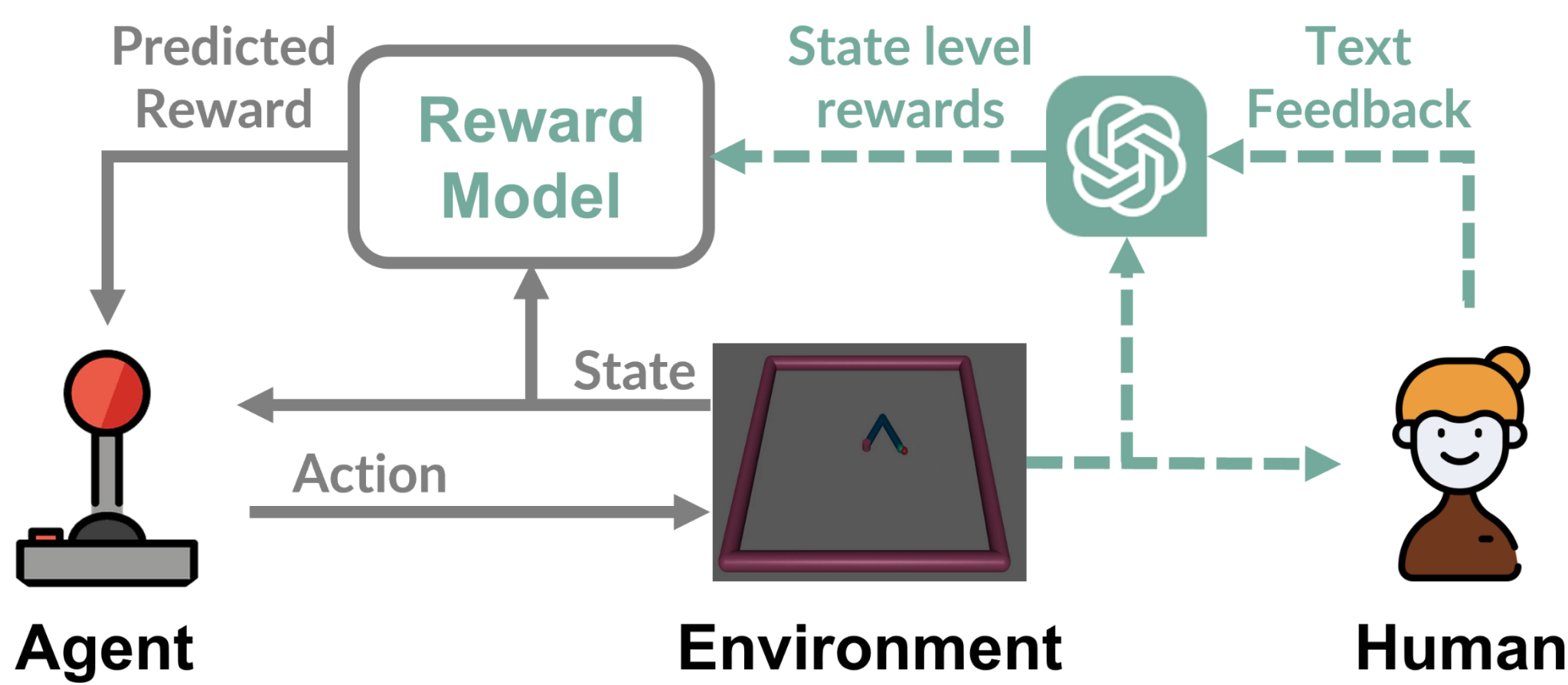


Reinforcement Learning from Human Text Feedback (RLHTF)

Belen Martin-Urcelay, Andreas Krause, Giorgia Ramponi



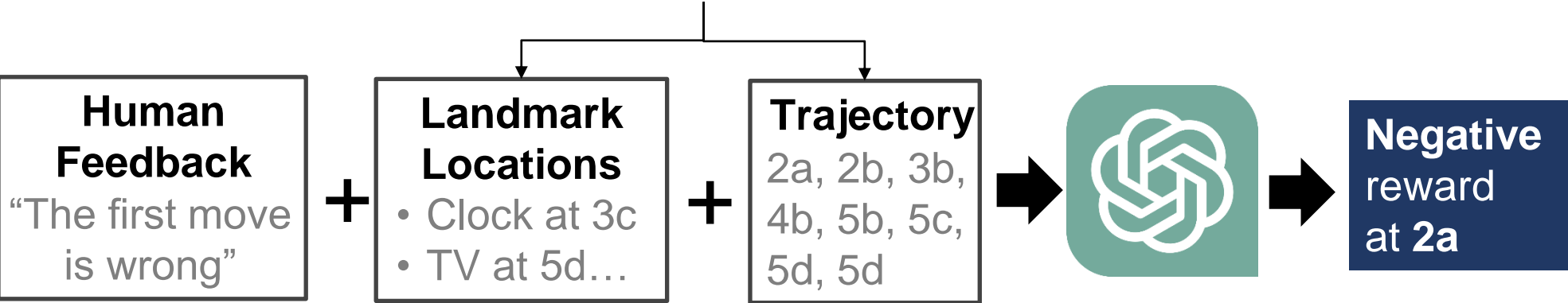
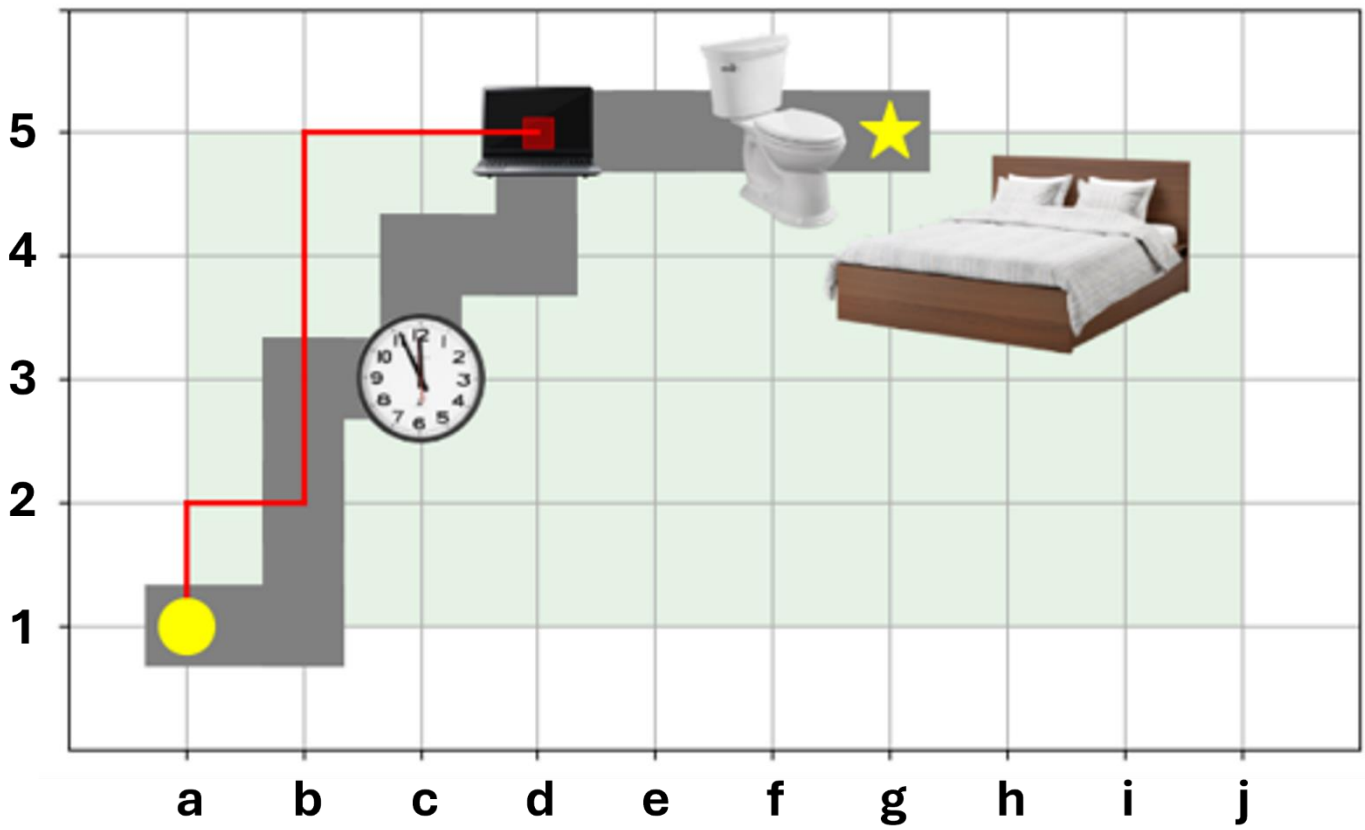
Framework



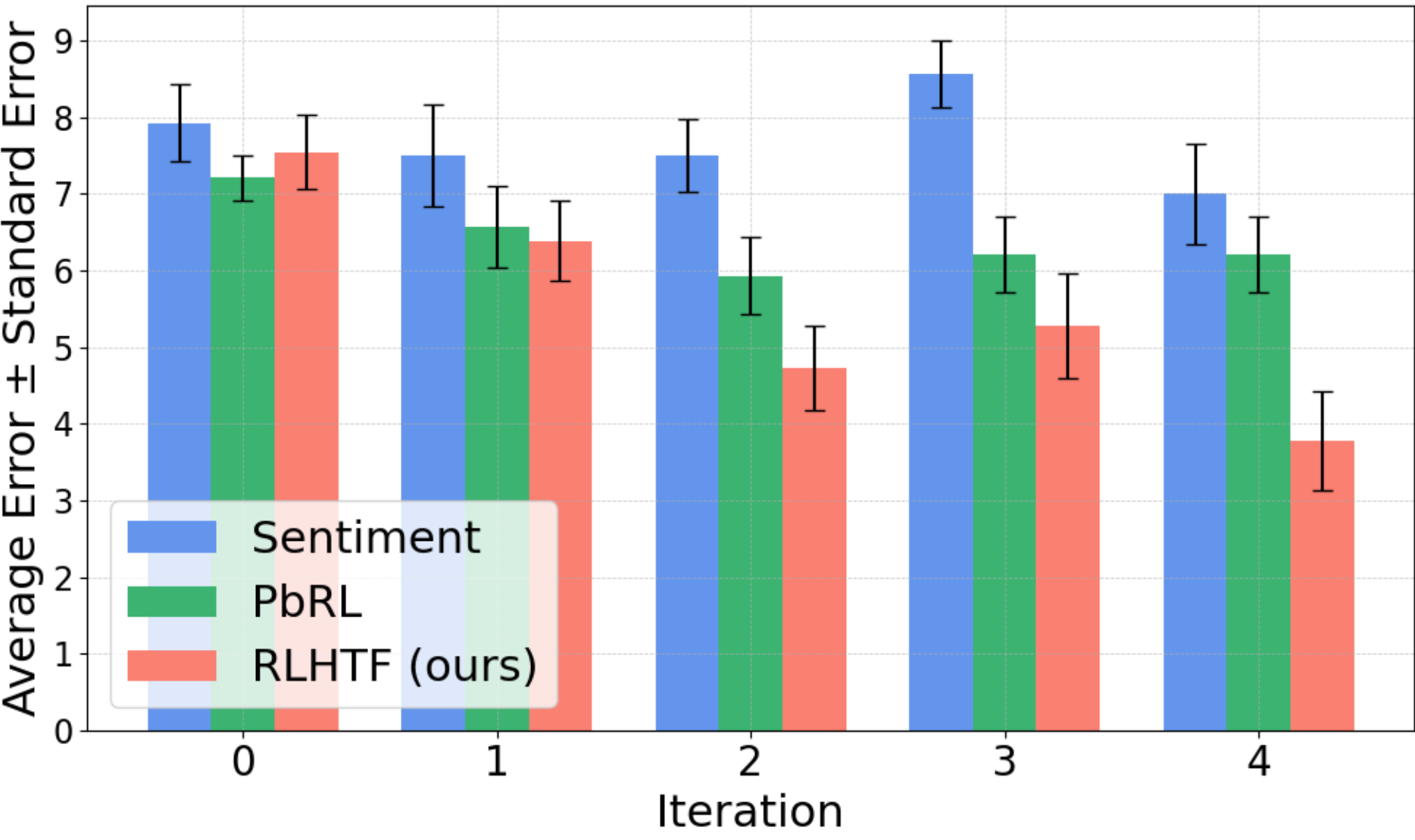
- An LLM translates human evaluations in the form of natural language into state level rewards.
- These labeled states are used to train a reward model.
- The agent is then trained with standard RL algorithms.

Experiments in Gridworld

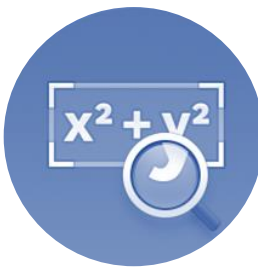
Goal: follow a specific path (in grey) from the start (yellow circle) to the end (yellow star). The agent's trajectory is in red.



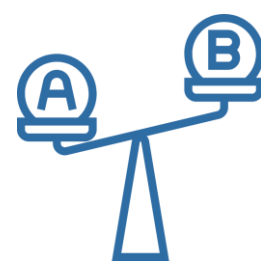
Performance Comparison



Motivation



Reward definition is challenging



Text has more information than rankings



LLMs are great at processing text

Use LLMs as a way of harnessing the information from human text to train a reward model efficiently in RLHF.

Baselines

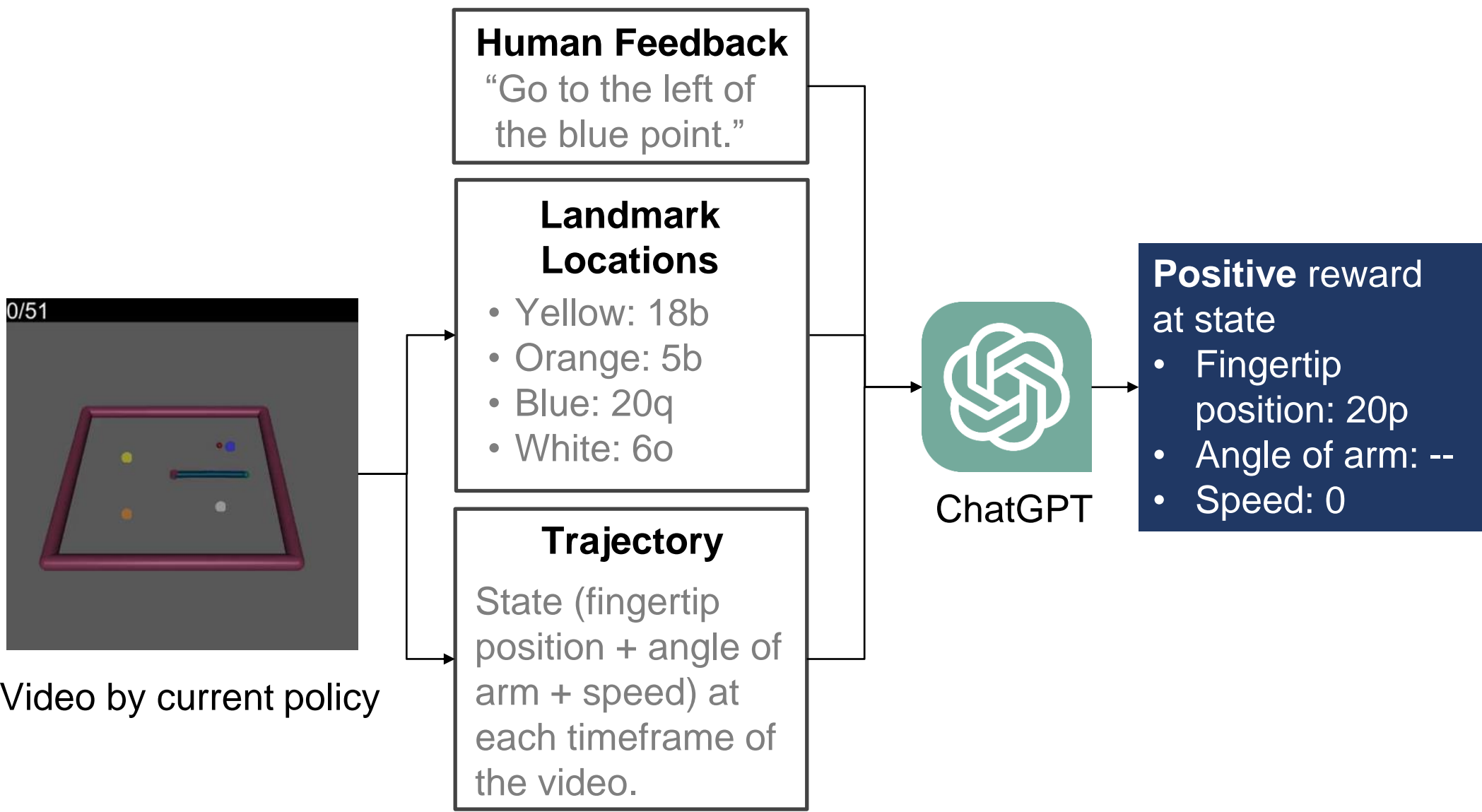
- **Sentiment:** We measure the sentiment of the human text feedback and apply it as a reward for all states in the trajectory [1].
- **PbRL:** We query human evaluators for pairwise comparisons between trajectories [2].
- **True:** The agent receives the true reward from the environment.

Algorithm

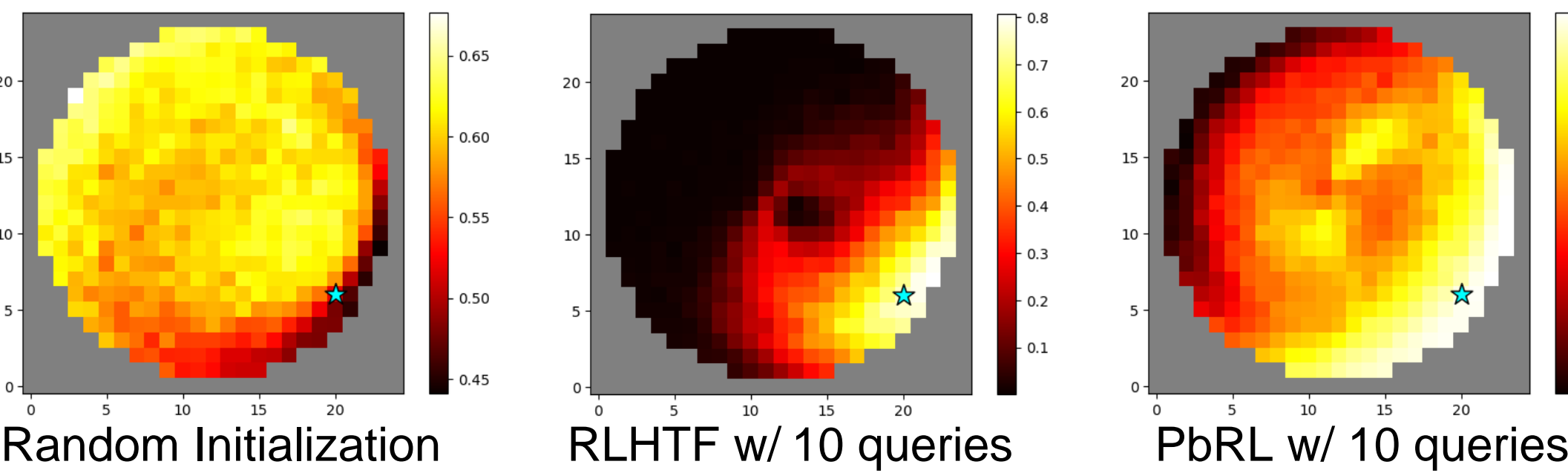
- 1: **Input:** number of human iterations N
- 2: Initialize Policy π_0 and reward model \hat{r}_0 .
- 3: **for** $i = 0$ to N **do**
- 4: Record trajectory following policy: $\mathbf{t}_i \leftarrow \pi_i$
- 5: Query human for feedback: $\mathbf{f}_i \leftarrow \mathbf{t}_i$
- 6: Translate feedback to state-reward pairs with LLM: $\mathbf{s}_i, \mathbf{r}_i \leftarrow \mathbf{f}_i, \mathbf{t}_i$
- 7: Update reward model: $\hat{r}_{i+1} \leftarrow \hat{r}_i, \{\mathbf{s}_t, \mathbf{r}_t\}_{t=0}^i$
- 8: Update policy: $\pi_{i+1} \leftarrow \pi_i, \hat{r}_{i+1}$
- 9: **end for**

Experiments in MuJoCo

Goal: Move a two-jointed robot arm to target (red point)

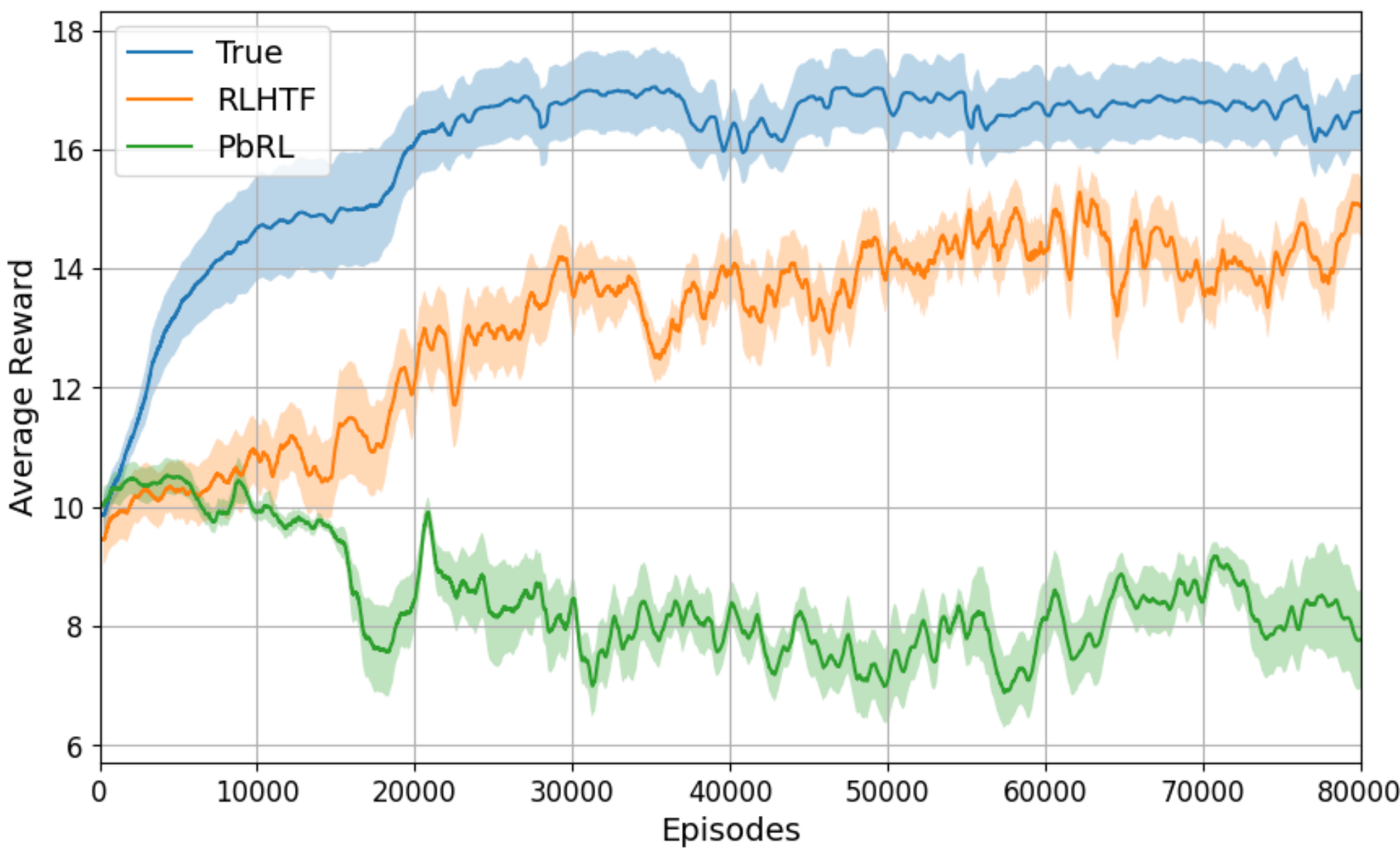


Reward Model Visualization



★: Target location □: Highest reward ■: Lowest reward

Performance Comparison



[1] Summers, T. R., Ho, M. K., Hawkins, R. D., Narasimhan, K., and Griffiths, T. L. Learning rewards from linguistic feedback. In Proc. of Conference on Artificial Intelligence
[2] Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. In Proc. of Advances in Neural Information Processing Systems.