

Boosting Feedback Efficiency of Interactive Reinforcement Learning by Adaptive Learning from Scores

Shukai Liu¹, Chenming Wu^{2†}, Ying Li³ and Liangjun Zhang²

Abstract—Interactive reinforcement learning has shown promise in learning complex robotic tasks. However, the process can be human-intensive due to the requirement of a large amount of interactive feedback. This paper presents a new method that uses scores provided by humans instead of pairwise preferences to improve the feedback efficiency of interactive reinforcement learning. Our key insight is that scores can yield significantly more data than pairwise preferences. Specifically, we require a teacher to interactively score the full trajectories of an agent to train a behavioral policy in a sparse reward environment. To avoid unstable scores given by humans negatively impacting the training process, we propose an adaptive learning scheme. This enables the learning paradigm to be insensitive to imperfect or unreliable scores. We extensively evaluate our method for robotic locomotion and manipulation tasks. The results show that the proposed method can efficiently learn near-optimal policies by adaptive learning from scores while requiring less feedback compared to pairwise preference learning methods. The source codes are publicly available at [https://github.com/SSKkai/Interactive-IRL](https://github.com/SSKkai/Interactive-Scoring-IRL).

I. INTRODUCTION

Deep Reinforcement Learning (DRL) has made remarkable progress in addressing robotic control tasks, such as legged robot locomotion [1] and robotic manipulation [2]. However, formulating an accurate reward function for a specific task can pose a significant challenge. Sparse reward functions are frequently employed for their simplicity, but their absence of reward signals can result in longer exploration and training time [3] and lower success rates [4]. In general, tasks with high complexity and high-dimensional continuous state-action spaces benefit from denser reward signals. Nonetheless, creating a reward function for autonomous agents necessitates domain expertise, which can be challenging for non-experts. Furthermore, hand-crafted rewards can be vulnerable to local optima or unexpected ways of achieving high numerical returns [5].

Inverse reinforcement learning (IRL) is a problem in which the reward function is inferred from expert demonstration trajectories [6]. This eliminates the need for tedious reward engineering and makes it possible to learn reward functions from expert demonstrations [7], [8]. However, IRL typically requires optimal demonstrations, while those demonstrated by people are often sub-optimal. To address this issue, a new method called Trajectory-ranked Reward Extrapolation

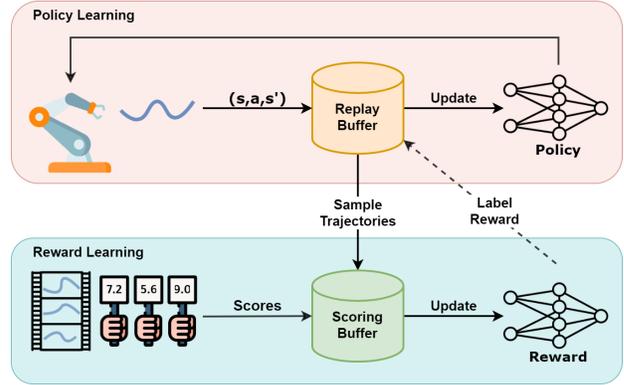


Fig. 1. Our proposed approach consists of two main components: policy learning and reward learning. The policy learning component updates itself through an off-policy RL algorithm to maximize the sum of predicted rewards generated by a reward network. Meanwhile, the reward learning component periodically samples trajectories, asks teachers to score them and optimizes its network.

(T-REX) was recently proposed. T-REX seeks to improve demonstrations by learning rewards from sequences of sub-optimal ranked demonstrations. It attempts to convert those rankings into a set of pairwise preferences to train the policy network [9]. However, T-REX still requires a large number of demonstrations to train the policy, which can be challenging in tasks where providing demonstrations at such a scale is difficult, even though optimality is not required.

In a recent study, PEBBLE [10] proposed an off-policy interactive RL algorithm to train a reward network and a policy network simultaneously from queried pairwise preferences. Teachers provide real-time pairwise feedback to supervise the learning process in the most efficient direction. However, this approach presents three major issues when providing feedback by assigning one-hot preference labels to two trajectories: 1) Two trajectories can be compared only when they have been paired together, making it difficult to gain a broader understanding of the relationship between individual and overall sampled trajectories. 2) Forcing the teacher to prioritize a better trajectory can sometimes become a burden and harm human-in-the-loop training. 3) To increase the number of training examples, partial trajectories (i.e., partitioning a full trajectory into segments) are used rather than full trajectories, which can make evaluating pairwise preferences more ambiguous for the teacher.

Our aim is to enhance feedback efficiency in interactive reinforcement learning. To achieve this, we suggest utilizing scores instead of pairwise preferences as the signals for interacting with RL agents. Moreover, we put forward

† Corresponding author

¹S. Liu is with the Department of Engineering Mathematics, University of Bristol, UK. lshukai22@gmail.com

²C. Wu and L. Zhang are with RAL, Baidu Research. [{wuchenming, liangjunzhang}@baidu.com">liangjunzhang}@baidu.com](mailto)

³Y. Li is with the School of Mechanical Engineering, Beijing Institute of Technology, Beijing, China. ying.li@bit.edu.cn

an adaptive learning scheme to make the training process smoother and more stable. This includes adaptive network optimization to smoothly update network parameters from score data and adaptive trajectory sampling to mine useful trajectories for teachers to evaluate, making our methodology less sensitive to imperfect or unreliable teacher inputs. By interchangeably giving feedback and adjusting the reward function, we continuously optimize both the reward and policy networks. Furthermore, we implement a scoring graphical user interface (GUI), showing the most relevant trajectories from the previously scored ones when a new trajectory is scored to support users in providing consistent scores. Teachers are allowed to amend and correct previous scores during the training process. An overview of our proposed framework is illustrated in Fig. 1. The main contributions of this paper are summarized as follows.

- 1) We develop an interactive RL method that enables the agent to learn both policy and reward simultaneously using a score-based approach. The RL agent proactively requests scores from teachers for complete trajectories, which results in requiring less feedback compared to pairwise-based methods.
- 2) We propose a method to tackle the problem of inaccuracies in scoring by introducing an adaptable learning approach that can withstand errors. Our proposed method also facilitates efficient learning of personalized and desired behaviors in situations where rewards are limited, based on the teachers' choices.

II. RELATED WORK

A. Inverse Reinforcement Learning

IRL allows the agent to better understand tasks and the environment, and learn an optimal policy using the reward via RL methods [11]. However, classic IRL frameworks [12], [8] assume that demonstrations are optimal and easy to obtain. Maximum entropy IRL [13], [14], [15] and Bayesian IRL [16], [17] are more robust to limited and stochastic suboptimality, but they cannot produce a policy better than the demonstration, thus their performance still highly relies on the quality of the demonstration. In [18], a generative model is learned from a large number of suboptimal demonstrations to produce noise-free trajectories. In [19], the reward function is formed as a linear combination of known features, and suboptimal demonstrations are utilized by learning rewards from trajectories labeled as success or failure. The method proposed in [20] is robust to a limited number of non-optimal demonstrations but still requires many expert demonstrations to discriminate the suboptimal ones. However, these methods can be challenging to apply in real-world scenarios where demonstrations are scarce, expensive, or suboptimal.

B. Learning from Evaluative Feedback

Evaluative feedback is a value given by a human teacher that rates the quality of the agent's behavior, which is easier for humans to provide compared to demonstrations. The TAMER framework [21] interprets evaluative feedback as a $Q^*(s, a)$ function of RL, and makes the agent act greedily

according to it. Meanwhile, the COACH framework [22] interprets human feedback as the advantage function $A^\pi(s, a)$ of policy gradient update. In the policy shaping framework [23], [24], evaluative feedback is considered an optimality label of the action. Providing evaluative feedback for entire trajectories can give each trajectory a global evaluation of its quality and, therefore, more accessible generalization. [25], [26] annotate each trajectory with a numeric score of a human teacher's global performance and leverage the IRL framework to learn a reward by minimizing the distance between human-provided and predicted scores. Although these methods show robustness to scoring errors, they cannot deal with suboptimal and high-dimensional IRL tasks.

Evaluative feedback is useful for handling non-optimality. CEILING [27] labels all state-action pairs with binary feedback evaluative feedback, then directly learns a Gaussian distributed policy by reinforcing the good ones while ignoring the bad ones. Similarly, [28] proposes IRLDC, which uses binary evaluative feedback to label each state-action pair. These methods still require a few demonstrations or correct feedback.

C. Preference-based Reinforcement Learning

[29] introduces the preference-based DRL framework, which can learn from pairwise preferences over the agent's current behaviors that the human teacher actively provides during training. This approach is on-policy that needs the human teacher to constantly provide preference feedback during the training process. Thus, [9] extends this framework and proposes T-REX, which learns the reward function from the pairwise preferences derived from a set of pre-collected ranked demonstrations, then applies the reward to RL for policy learning. T-REX allows the learned reward to extrapolate beyond the demonstrations and achieve better-than-demonstrator performance. D-REX [30] and SSRR [31] extend the learning-from-ranking framework by automatically generating ranked trajectories via noise injection. However, the need for demonstrations still exists.

Myers et al. [32] proposed a robot learning method that can learn multimodal rewards from multiple active ranking queries by multiple experts. PEBBLE [10] presented an interactive preference learning method that enables users to give preference feedback directly on the behavior of the RL agent, thus eliminating the need for demonstrations. PEBBLE introduces the off-policy learning framework to reuse data and follows the feedback form of pairwise preferences between partial trajectories for the sample efficiency as in [29]. To improve feedback efficiency, [33] investigates the query selection and policy initialization. [34] presents an exploration method to collect more diverse experiences. [35] introduces this learning scheme for socially aware robot navigation and reduces the amount of preference feedback from humans by collecting expert demonstrations. [36] further increases the feedback efficiency by inferring pseudo-labels on a large number of unlabeled samples with data augmentation, while we annotate global scores to the agent's past experience. In our work, we annotate global scores to

the agent’s past experiences and demonstrate that this scoring feedback scheme can substantially reduce the amount of required feedback and better fit the off-policy framework.

III. METHODOLOGY

Our proposed framework can be broken down into two processes. First, the RL agent interacts with the environment to create new trajectories to be scored. Second, an off-policy DRL algorithm is applied to update the agent’s policy π_ψ in order to maximize the expectation of the predicted reward generated by \hat{r}_θ . Additionally, the teacher reviews the sampled trajectories through video replay and scores them at a frequency of f during the RL training process. The scored trajectories (τ, s) are stored in the scoring buffer \mathcal{D} to update the reward network. The agent then deduces the teacher’s preference from the score difference and updates the reward network accordingly. The updated reward network guides the agent to generate better trajectories. Scoring these trajectories can lead to a more comprehensive reward network, and the agent learns the policy and reward simultaneously.

A. Adaptive Learning from Scores

The reward function is trained as a neural network, which we refer to as \hat{r}_θ . Users can choose either states or state-action pairs as input. Our approach utilizes two replay buffers: one for the RL part to store the state-action transitions, and the other for reward learning to store the trajectories and their scores. For policy learning, the only difference between our method and vanilla RL algorithms is that the rewards are produced by the reward network. As a result, our approach can be applied to most off-policy RL algorithms while maintaining their core functionality.

Note that the reward function is dynamically updated during training, which can cause inconsistency in off-policy RL since previously generated rewards may not match the latest reward functions. To address this issue, we adopt the approach proposed in PEBBLE [10] and relabel the replay buffer each time we update the reward. Storing all scored trajectories in the scoring buffer allows for off-policy learning in reward learning. This allows newly scored trajectories to be compared with previously scored trajectories, which significantly improves the utilization rate of human feedback. Moreover, the scoring buffer allows the teacher to access previous scores during training and correct them if they change their minds, making reward learning more robust.

Adaptive Network Optimizing. Given a set of trajectories $\tau_1, \tau_2, \dots, \tau_m$ and their corresponding scores s_1, s_2, \dots, s_m , our goal is to parameterize a reward function \hat{r}_θ to infer the underlying reward from scores and output the reward value that matches user’s evaluation standard, such that $\sum_{\tau_i} \hat{r}_\theta < \sum_{\tau_j} \hat{r}_\theta$ when $s_i < s_j$. This problem can be regarded as a learning-to-rank problem [37], to optimize the following equation:

$$\min U(\text{sorted}(\hat{r}_\theta(\tau_1), \dots, \hat{r}_\theta(\tau_{|\mathcal{D}|})), \mathbf{y}) \quad (1)$$

where \mathbf{y} is the ground-truth index list and U is a binary function to evaluate if the ranked position is equivalent to

the ground-truth position y_i in \mathbf{y} . In our work, we decide to solve this problem using the idea of pairwise learning because our major goal is to train a reward function, instead of a ranker that generates a descent permutation. The user’s preference over any pair of two trajectories is described by a distribution μ , which can be derived by comparing their scores, e.g., $\mu = 1$ if $s_i < s_j$. By following the Bradley-Terry and Luce-Shephard models of preferences [38], [39], a preference predictor using the reward function \hat{r}_θ can be modeled as a softmax-normalized distribution as follows.

$$P(\tau_i \prec \tau_j) = \frac{e^{\sum_{\tau_j} \hat{r}_\theta}}{e^{\sum_{\tau_i} \hat{r}_\theta} + e^{\sum_{\tau_j} \hat{r}_\theta}} \quad (2)$$

where $\tau_i \prec \tau_j$ denotes the trajectory τ_j is more preferred to the trajectory τ_i . This equation demonstrates that the probability of preferring one trajectory to another is exponentially related to the predicted return of each trajectory. Thus, the parameterized reward function \hat{r}_θ can be learned via minimize the cross entropy loss between the predicted preference and the user’s true preference as follows.

$$\mathcal{L} = - \sum_{(\tau_i, \tau_j, \mu) \in \mathcal{D}} [\mu \log P(\tau_i \prec \tau_j) + (1 - \mu) \log P(\tau_j \prec \tau_i)] \quad (3)$$

The preference distribution μ is usually a one-hot encoded label. Although this can learn reward effectively on correct labels, it may suffer from poor performance when there are wrong labels in the database \mathcal{D} . Unfortunately, it is nearly impossible for a human user to score a large number of trajectories perfectly. To strengthen the robustness against scoring error, we use the label smoothing method [40] to convert the hard label μ to soft label $\tilde{\mu}$ using $\tilde{\mu} = (1 - \alpha)\mu + \alpha/K$, where $\alpha \in [0, 1]$ is a constant factor that indicates the smoothing strength for one-hot label and K denotes the number of labels, in our case, $K = 2$. However, in our setting, using a constant smoothing strength for all pairwise labels may not be ideal because it ignores the relative relationship implicit in the score differences. It is intuitive that for a trajectory pair, the larger the scores differ, the more confident that one trajectory is better than the other. Thus, to better exploit the information from human scores, we make the smoothing strength adaptive to the score differences by $\alpha = 1/(|s_i - s_j| + \lambda)^2$ where $\lambda > 1$ is a hyperparameter, and we set it as 2 in all our experiments. The adaptive α makes the label $\tilde{\mu}$ closer to 0 or 1 when pairwise trajectories significantly differ in the score and approach 0.5 when the scores are similar. Hence, the soft label $\tilde{\mu}$ is computed as

$$\tilde{\mu} = (1 - \frac{1}{(|s_i - s_j| + \lambda)^2})\mu + \frac{1}{K(|s_i - s_j| + \lambda)^2} \quad (4)$$

Adaptive Trajectory Sampling. The process of learning through rewards includes two sampling procedures. The first one involves gathering newly created trajectories from the RL agent, which are then evaluated and given scores by

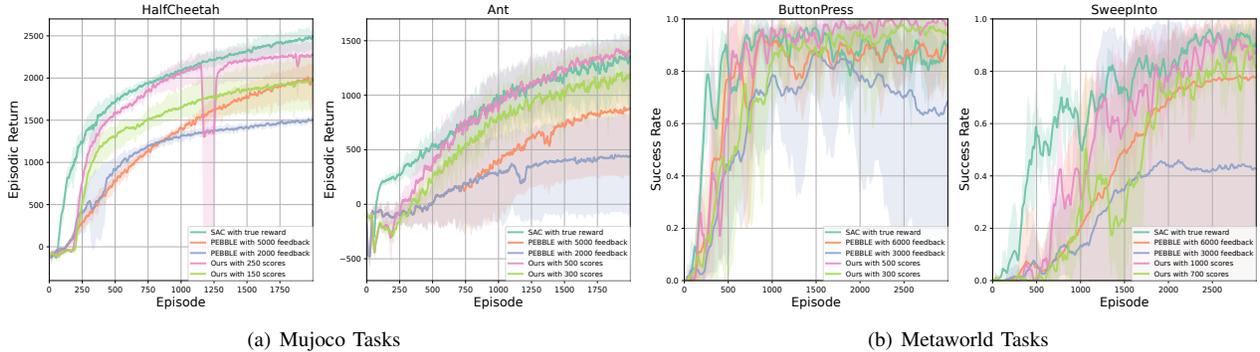


Fig. 2. The graphs display the learning progress for SAC, PEBBLE, and our method in two different simulations: HalfCheetah and Ant for movement tasks, and ButtonPress and SweepInto for robotics tasks. The data is measured by the success rate for the Metaworld environment and ground true return for Mujoco. The solid lines indicate the mean, while the shaded regions show the minimum and maximum range across three runs for all figures.

teachers. The second procedure involves selecting a batch of trajectory pairs from the scoring buffer, which stores all scored trajectories.

To improve the training of our RL agent, we ask the user to rate the newly generated trajectories. These ratings are stored in the scoring buffer called \mathcal{D} as scored trajectories (τ, s) . However, asking for scores for all trajectories can be overwhelming. Therefore, we aim to choose the most informative scoring queries. This way, even if only a few newly generated trajectories are scored each time, they are enough to train the appropriate reward. We employ the k -means clustering algorithm to automatically select trajectories with high variance in performance, which are then approximated by the predicted rewards. To select k trajectories from a set of newly generated ones for evaluation, we use the reward network to compute the episodic return of each trajectory. Next, we run k -means clustering on these returns and choose the k trajectories whose returns are closest to each k centroid.

For n scored trajectories (τ, s) in the scoring buffer \mathcal{D} , we notice that not all the trajectory pairs of them can lead to effective reward update. To address this, we explore methods for sampling from the scoring buffer \mathcal{D} . An off-the-shelf sampling method is the entropy-based sampling adopted in PEBBLE [10], which randomly samples a large batch of trajectories and seeks to maximize the entropy. However, we notice that when a trajectory with a higher score is sampled into \mathcal{D} , it should be compared more broadly with other trajectories. This allows the reward equation to learn what behaviors lead to higher scores. Unfortunately, entropy-based sampling cannot provide this capability. Inspired by the prioritized experience replay (PER) methods [41], we propose an alternative sampling methodology: either randomly selecting one trajectory in a pair or choosing based on a probability that increases with its score. The probability of each scored trajectory is computed according to its score as $P(i) = \frac{s_i^\beta}{\sum_n s_n^\beta}$, where β is a hyperparameter that determines how much prioritization is assigned to a high-scored trajectory. And n is the total number of trajectories in the scoring buffer. The comparison between our sampling method and entropy-based sampling is shown in the experimental section.

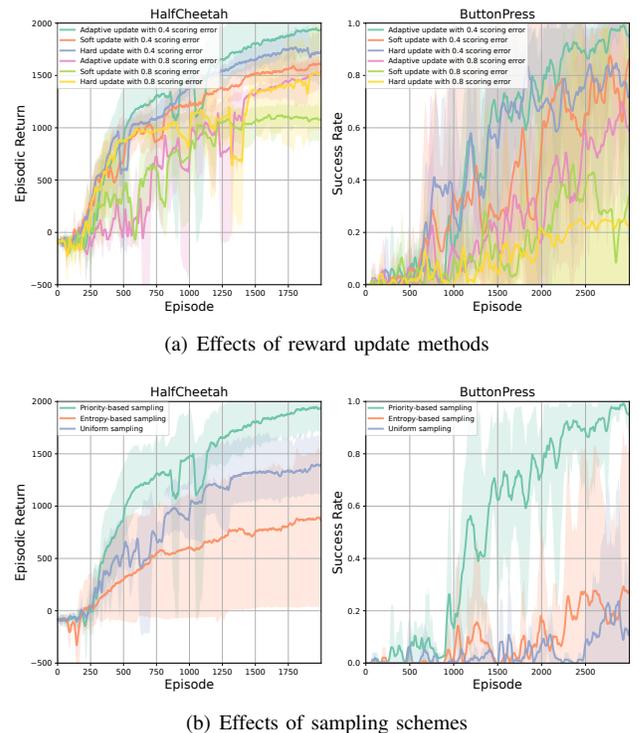


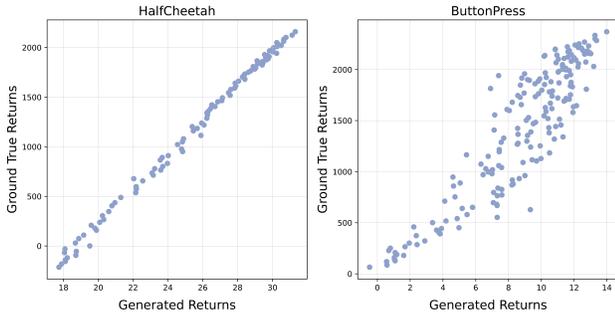
Fig. 3. Ablation study on HalfCheetah and ButtonPress. (a) The performance comparisons of different update reward methods under the scoring noise of 0.4 and 0.8. (b) Effects of sampling schemes to select training batch for reward update under the scoring noise of 0.4. We set a maximum of 400 scores for HalfCheetah and 1000 scores for ButtonPress.

IV. EXPERIMENTS

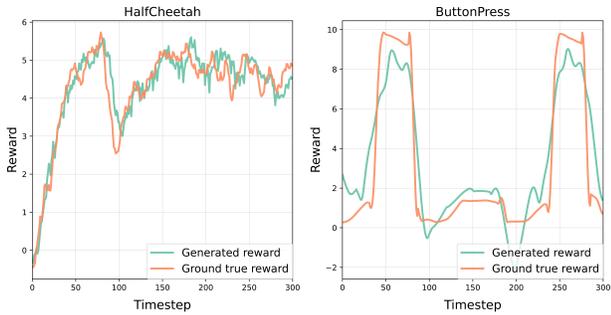
A. Experiment Setups

We compare our approach to previous methods to verify if our approach can achieve similar performance with less feedback. In Sec. IV-B, we conduct ablation studies to investigate the influence of adaptive reward updates on the robustness of scoring errors, and to assess how various sampling methods affect performance. In Sec. IV-D, we analyze the learned reward and the agent’s behavioral pattern to determine if our approach can accurately extrapolate the user’s preferences and underlying intent. Finally, we conduct real human experiments in Sec. IV-E.

We evaluate our proposed method on several continuous



(a) Analysis of episodic returns



(b) Analysis of reward signals within a episode

Fig. 4. (a) The episodic returns of learned reward and ground true reward (b) The learned reward signals and the ground true reward signals within a single episode along the timesteps.

robotic tasks in simulation, including locomotion tasks Ant and HalfCheetah in Mujoco simulator [42] with OpenAI Gym [43], and robotic manipulation tasks in Metaworld environment [44], namely PushButton and SweepInto. For locomotion tasks, we use the episode return as the evaluation metric. For manipulation tasks, we use the task success rate of the last 100 episodes as the evaluation metric. We train 2,000 episodes for Mujoco locomotion tasks and 3,000 episodes for Metaworld robotic manipulation tasks each run. The episode is 300 steps long, with the exception of the SweepInto task, which is 250 steps long to reduce the proportion of task-goal-unrelated steps in the episode.

We use the state-of-the-art off-policy DRL algorithm SAC to learn the behavioral policy [45]. However, the agent can only receive the reward generated by our learned reward function. To model our reward function, we use a single deep neural network consisting of 3 fully connected layers of 256 units with leaky ReLUs. We train the reward network from scores using the Adam optimizer with a learning rate of 10^{-3} and a batch size of 128. We use a standard scoring range of 0 to 10 across all experiments. To evaluate our approach quantitatively, we make the agent learn tasks only by intermittently getting its past experience scored by a scripted teacher. By linearly mapping the episode returns to the scoring range, the scripted teacher provides scores. We use a two-stage scoring frequency: at the beginning of training, we use a faster scoring frequency, scoring 5 trajectories every 10 episodes. When the agent’s performance reaches approximately a quarter of the maximum episode returns, we switch to a slower scoring frequency, scoring 10

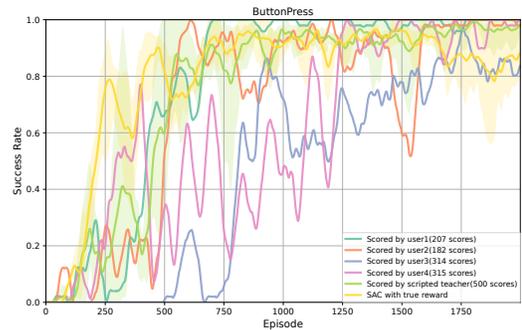


Fig. 5. The performance comparisons of original SAC and our approach trained by real human users and scripted teacher.

trajectories every 100 episodes. For all experiments, we set the β for adaptive sampling to 3.

B. Results

To examine the effectiveness of our approach, we compare it to the original SAC algorithm training with the same ground true reward as we used for the scripted teacher. We use the same hyperparameters for both trainings. We also compare to the state-of-art preference learning algorithm PEBBLE [10]. We use the exact same values of hyperparameters for PEBBLE as the Equal SimTeacher setting reported in [33] and the corresponding open-source code repository.

Fig. 2 shows the learning curves of our approach and PEBBLE with different numbers of teacher preference feedback in comparison to SAC with true reward. Note that our approach employs a different type of feedback than PEBBLE. In this experiment setting, PEBBLE assigns a preference label to a pair of 50-step partial trajectories for single teacher feedback, whereas our approach assigns a global score to an entire episode with 300 steps. As a result, we give PEBBLE an advantage by providing more than three times as much feedback as ours. We can see that our approach achieves the same or higher level of performance than PEBBLE, which is given more feedback, in all tasks. In comparison to the SAC with ground true reward, our approach requires more training time to converge because it must learn the reward from scratch at the beginning of training, but it can match the performance after convergence in all tasks using only a small number of trajectory scores from the teacher. The results show that our approach can learn robot behavioral policies effectively in sparse reward environments with teachers’ scores.

C. Ablation Study

1) *Robustness to Scoring Errors*: The preceding experiments assume access to the perfect correct scores generated by the ground true reward. However, in practice, it is impossible for a human teacher to score hundreds of trajectories accurately: users can give vague and approximate scores for trajectories with similar performance. Thus, we examine the robustness to scoring errors and low scoring precision of our approach by comparing the performance of the noisy scores when using hard reward update, soft reward update by label smoothing, and adaptive reward update.

We simulate the real human teacher by adding noise randomly generated by Gaussian distribution to the scores given by the scripted teacher as $s' = \mathcal{N}(s, \sigma_{noise}^2)$. We adopt a minimal step of 0.5 for these noise-infused scores, such as 3.0 and 7.5. This permits teachers to give scores to trajectories that perform similarly. We tested our method with $\sigma_{noise}^2 = 0.4$ and $\sigma_{noise}^2 = 0.8$, and use Kendall’s τ_B coefficient to further measure the rank correlation between the noisy scores and the perfect correct scores. This coefficient is calculated by $\tau_B = (P - Q) / \sqrt{(P + Q + T)(P + Q + U)}$, where P is the number of concordant pairs, Q is the number of discordant pairs, T is the number of ties only in one group of data, and U is the number of ties only in another, τ_B will be high and close to 1 if two variables have similar rank. We found that in our experiment, $\sigma_{noise}^2 = 0.4$ corresponds to $\tau_B \approx 0.8$. The higher noise $\sigma_{noise}^2 = 0.8$ leads to a lower correlation level $\tau_B \approx 0.65$.

The results of training with imperfect scores on the HalfCheetah and ButtonPress tasks are shown in Fig. 3(a), where the smoothing strength is $\alpha = 0.05$ for the original label smoothing, $\alpha' = 2$ for the adaptive reward update. Despite a slight decrease in performance compared to perfect scoring, the adaptive update method performs better compared to the other two methods in both tasks when scoring errors $\sigma_{noise}^2 = 0.4$. With $\sigma_{noise}^2 = 0.8$, the adaptive reward update method surpasses others in the ButtonPress task. However, in the relatively simple task, HalfCheetah did not gain extra advantages over the hard reward update. Overall, our proposed adaptive reward update method delivers the strongest performance and shows strong robustness to high-scoring errors. We also investigate the effects of different sampling methods to select scored trajectory pairs for reward updates. We used $\sigma_{noise}^2 = 0.4$ to simulate the real scoring scenario. Fig. 3(b) shows the learning curves of our approach on the HalfCheetah and ButtonPress tasks under three different sampling schemes: uniform sampling, entropy-based sampling, and priority-based sampling. We can see that the priority-based sampling method significantly outperforms other sampling methods. Although entropy-based sampling performs well with a perfect feedback teacher as suggested in [33], it cannot handle the noisy-scoring scenario well.

D. Reward Extrapolation

1) *Reward Analysis:* We compare the learned reward function to the ground truth rewards to assess the quality of the learned reward function. We run SAC with true reward to collect trajectories with a variety of performance qualities, and then we compare the episodic ground truth returns to the returns generated by the learned reward function. Fig. 4(a) shows the reward function learned by our approach on 250 scores and 500 scores in HalfCheetah and ButtonPress respectively. We can see that the learned reward function has a strong correlation with the true reward on the ground. It should be noted that the learned and true rewards have very different scales, but this difference had no effect on policy learning performance. We further investigate the rewards by looking into the reward functions within an episode

at different timesteps. We generate a set of suboptimal trajectories with high and low reward ranges in one episode. The results are shown in Fig. 4(b). We manually normalize the learned reward outputs to have the same scale as the true rewards by multiplying a coefficient. We can see that the learned rewards are well-aligned with the ground truth rewards.

2) *Customized Behavior:* One goal of our approach is to enable users to train customized policies through scoring. We demonstrate this in the RL Bench [46] simulation task PushButton, which requires a Franka Emika Panda robot arm to push a button on a table. We model two scripted teachers to score trajectories with different preferences as follows: (1) teacher 1: robot first moves to the top of the button, then pushes with the gripper tip while remaining vertical, (2) teacher 2: robot moves its gripper parallel to the table and presses the button with its side. For the trained agent’s policies please refer to the supplementary video. The result demonstrates that our method can infer users’ underlying intent and complete tasks in accordance with the user’s preferences.

E. Real Human Experiment

We conduct experiments with real human users to test our approach. We create a graphical user interface (GUI) and test it with two users on the MetaWorld ButtonPress environment. The GUI displays four previously scored trajectories and their scores as references to help users score new trajectories consistently. We select two scored trajectories with the closest predicted returns to the current trajectory and two references that are most similar to the current trajectory in Cartesian space, measured by dynamic time warp (DTW) [47]. Users are allowed to revise the scores of the reference trajectories as needed, and they could skip scoring if they find it difficult. We follow the two-stage scoring frequency outlined in Sec. IV-A, starting with a faster scoring frequency and allowing users to switch to a lower frequency mode based on their performance. Fig. 5 shows the learning curve of the four users compared to learning by SAC with true reward and learning by our approach with a scripted teacher. Our approach shows that a good behavior policy could be trained with only about three hundred scores. For more information on using the scoring interface, please refer to the supplementary video.

V. CONCLUSION

We propose an algorithm for interactive RL that uses scores from a teacher to learn both a policy and reward function. This eliminates the need for human demonstrations and maximizes the use of user feedback, reducing the amount of required feedback. Our experiments show that even with a small number of human scores, our method can train robotic locomotion and manipulation tasks to near-optimal levels. With this method, we can map global behavior evaluations to rewards for only states or state-action pairs, allowing us to learn optimal policies in environments where rewards cannot be observed.

REFERENCES

- [1] C. Yang, K. Yuan, Q. Zhu, W. Yu, and Z. Li, "Multi-expert learning of adaptive legged locomotion," *Science Robotics*, vol. 5, no. 49, p. eabb2174, 2020.
- [2] R. Liu, F. Nageotte, P. Zanne, M. de Mathelin, and B. Dresp-Langley, "Deep reinforcement learning for the control of robotic manipulation: a focussed mini-review," *Robotics*, vol. 10, no. 1, p. 22, 2021.
- [3] M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, "Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards," *arXiv preprint arXiv:1707.08817*, 2017.
- [4] A. Mohtasib, G. Neumann, and H. Cuayáhuitl, "A study on dense and sparse (visual) rewards in robot policy learning," in *Annual Conference Towards Autonomous Robotic Systems*. Springer, 2021, pp. 3–13.
- [5] D. Hadfield-Menell, S. Milli, P. Abbeel, S. J. Russell, and A. Dragan, "Inverse reward design," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [6] S. Arora and P. Doshi, "A survey of inverse reinforcement learning: Challenges, methods and progress," *Artificial Intelligence*, vol. 297, p. 103500, 2021.
- [7] A. Y. Ng, S. J. Russell *et al.*, "Algorithms for inverse reinforcement learning," in *Int. Conf. on Machine Learning*, vol. 1, 2000, p. 2.
- [8] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Int. Conf. on Machine Learning*, 2004, p. 1.
- [9] D. Brown, W. Goo, P. Nagarajan, and S. Niekum, "Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations," in *Int. Conf. on Machine Learning*. PMLR, 2019, pp. 783–792.
- [10] K. Lee, L. Smith, and P. Abbeel, "Pebble: Feedback-efficient interactive reinforcement learning via relabeling experience and unsupervised pre-training," *arXiv preprint arXiv:2106.05091*, 2021.
- [11] R. S. Sutton, A. G. Barto *et al.*, "Introduction to reinforcement learning," 1998.
- [12] A. Y. Ng, S. J. Russell *et al.*, "Algorithms for inverse reinforcement learning," in *Int. Conf. on Machine Learning*, vol. 1, 2000, p. 2.
- [13] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey *et al.*, "Maximum entropy inverse reinforcement learning," in *Proc. of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 8. Chicago, IL, USA, 2008, pp. 1433–1438.
- [14] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum entropy deep inverse reinforcement learning," *arXiv preprint arXiv:1507.04888*, 2015.
- [15] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization," in *Int. Conf. on Machine Learning*, 2016, pp. 49–58.
- [16] D. Ramachandran and E. Amir, "Bayesian inverse reinforcement learning," in *Int. Conf. on Artificial Intelligence*, vol. 7, 2007, pp. 2586–2591.
- [17] B. Okal and K. O. Arras, "Learning socially normative robot navigation behaviors with bayesian inverse reinforcement learning," in *Int. Conf. on Robotics & Automation*. IEEE, 2016, pp. 2889–2895.
- [18] A. Coates, P. Abbeel, and A. Y. Ng, "Learning for control from multiple demonstrations," in *Int. Conf. on Machine Learning*, 2008, pp. 144–151.
- [19] K. Shiarlis, J. Messias, and S. Whiteson, "Inverse reinforcement learning from failure," 2016.
- [20] S. Choi, K. Lee, and S. Oh, "Robust learning from demonstrations with mixed qualities using leveraged gaussian processes," *IEEE Transactions on Robotics*, vol. 35, no. 3, pp. 564–576, 2019.
- [21] W. B. Knox and P. Stone, "Interactively shaping agents via human reinforcement: The tamer framework," in *Proceedings of the fifth international conference on Knowledge capture*, 2009, pp. 9–16.
- [22] J. MacGlashan, M. K. Ho, R. Loftin, B. Peng, G. Wang, D. L. Roberts, M. E. Taylor, and M. L. Littman, "Interactive learning from policy-dependent human feedback," in *Int. Conf. on Machine Learning*. PMLR, 2017, pp. 2285–2294.
- [23] S. Griffith, K. Subramanian, J. Scholz, C. L. Isbell, and A. L. Thomaz, "Policy shaping: Integrating human feedback with reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 26, 2013.
- [24] T. Cederborg, I. Grover, C. L. Isbell, and A. L. Thomaz, "Policy shaping with human teachers," in *Int. Conf. on Artificial Intelligence*, 2015.
- [25] L. El Asri, B. Piot, M. Geist, R. Laroche, and O. Pietquin, "Score-based inverse reinforcement learning," in *International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2016)*, 2016.
- [26] B. Burchfiel, C. Tomasi, and R. Parr, "Distance minimization for reward learning from scored trajectories," in *Proc. of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 30, no. 1, 2016.
- [27] E. Chisari, T. Welschehold, J. Boedecker, W. Burgard, and A. Valada, "Correct me if i am wrong: Interactive learning for robotic manipulation," *IEEE Robotics and Automation Letters*, 2022.
- [28] N. Mourad, A. Ezzeddine, B. Nadjar Araabi, and M. Nili Ahmadabadi, "Learning from demonstrations and human evaluative feedbacks: Handling sparsity and imperfection using inverse reinforcement learning approach," *Journal of Robotics*, vol. 2020, 2020.
- [29] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [30] D. S. Brown, W. Goo, and S. Niekum, "Better-than-demonstrator imitation learning via automatically-ranked demonstrations," in *Conference on robot learning*, 2020, pp. 330–359.
- [31] L. Chen, R. Paleja, and M. Gombolay, "Learning from suboptimal demonstration via self-supervised reward regression," *arXiv preprint arXiv:2010.11723*, 2020.
- [32] V. Myers, E. Biyik, N. Anari, and D. Sadigh, "Learning multimodal rewards from rankings," in *Conference on Robot Learning*. PMLR, 2022, pp. 342–352.
- [33] K. Lee, L. Smith, A. Dragan, and P. Abbeel, "B-pref: Benchmarking preference-based reinforcement learning," *arXiv preprint arXiv:2111.03026*, 2021.
- [34] X. Liang, K. Shu, K. Lee, and P. Abbeel, "Reward uncertainty for exploration in preference-based reinforcement learning," *arXiv preprint arXiv:2205.12401*, 2022.
- [35] R. Wang, W. Wang, and B.-C. Min, "Feedback-efficient active preference learning for socially aware robot navigation," in *Int. Conf. on Intelligent Robots and Systems*. IEEE, 2022, pp. 11 336–11 343.
- [36] J. Park, Y. Seo, J. Shin, H. Lee, P. Abbeel, and K. Lee, "Surf: Semi-supervised reward learning with data augmentation for feedback-efficient preference-based reinforcement learning," *arXiv preprint arXiv:2203.10050*, 2022.
- [37] Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li, "Learning to rank: from pairwise approach to listwise approach," in *Int. Conf. on Machine Learning*, 2007, pp. 129–136.
- [38] R. A. Bradley and M. E. Terry, "Rank analysis of incomplete block designs: I. the method of paired comparisons," *Biometrika*, vol. 39, no. 3/4, pp. 324–345, 1952.
- [39] R. D. Luce, *Individual choice behavior: A theoretical analysis*. Courier Corporation, 2012.
- [40] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.
- [41] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," *arXiv preprint arXiv:1511.05952*, 2015.
- [42] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *Int. Conf. on Intelligent Robots and Systems*, 2012, pp. 5026–5033.
- [43] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016.
- [44] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine, "Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning," in *Conference on Robot Learning (CoRL)*, 2019.
- [45] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Int. Conf. on Machine Learning*, 2018, pp. 1861–1870.
- [46] S. James, Z. Ma, D. Rovick Arrojo, and A. J. Davison, "Rlbench: The robot learning benchmark & learning environment," *IEEE Robotics and Automation Letters*, 2020.
- [47] D. J. Berndt and J. Clifford, "Using dynamic time warping to find patterns in time series," in *KDD workshop*, vol. 10, no. 16, 1994, pp. 359–370.