

---

# Tell my why: Training preferences-based RL with human preferences and step-level explanations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Human-in-the-loop reinforcement learning (HRL) allows the training of agents  
2 through various interfaces, even for non-expert humans. Recently, preference-  
3 based methods (PBRL), where the human has to give his preference over two  
4 trajectories, increased in popularity since they allow training in domains where  
5 more direct feedback is hard to formulate. However, the current PBRL methods  
6 have limitations and do not provide humans with an expressive interface for giving  
7 feedback. With this work, we propose a new preference-based learning method  
8 that provides humans with a more expressive interface to provide their preference  
9 over trajectories and a factual explanation (or annotation of why they have this  
10 preference). These explanations allow the human to explain what parts of the  
11 trajectory are most relevant for the preference. We allow the expression of the  
12 explanations over individual trajectory steps. We evaluate our method in various  
13 simulations using a simulated human oracle (with realistic restrictions), and our  
14 results show that our extended feedback can improve the speed of learning. Code  
15 & data: [github.com/under-rewiev](https://github.com/under-rewiev)

## 16 1 Introduction

17 To allow for more flexible deployments of agents in all areas of human life, it is highly desirable that  
18 these can be quickly and easily adapted to different needs. While there are many different methods  
19 for training agents, Human-in-the-Loop Reinforcement Learning puts the responsibility on the human  
20 who gives feedback (or other methods like advice or demonstrations) from which the robots can learn  
21 (Amershi et al. [2014]). This has the advantage that there is no need to program hard-coded actions  
22 or extensive modeling for planning or designing reward functions, which can be a cumbersome  
23 and brittle process in itself (Booth et al. [2023]). Pushing the responsibility directly to the human  
24 in the training process alleviates some of these problems while opening the learning process up to  
25 more non-technical users (or trainers). Naturally, there are multiple options for how a human can  
26 direct their feedback toward the agent. In this work, we want to focus on one of the more prominent  
27 methods of teaching agents by giving preferences over two trajectories, often called preferences-based  
28 reinforcement learning (PBRL). While Human-in-the-loop Reinforcement Learning (HRL) recently  
29 gained more interest from the research community, most research is focused on purely technical  
30 aspects with the same interface. We argue that extending that interface and giving humans more  
31 options to control learning is important.

32 Our work is also motivated by the concept of "scaffolding" (van de Pol et al. [2010]), which human  
33 teachers often use when teaching tasks to students. While scaffolding contains many techniques, like  
34 breaking the task into sub-parts, using verbal cues, or relating the content to other knowledge, we  
35 focus on a single technique in this work. Teachers often highlight essential features of the tasks (or  
36 solutions) to the student. This highlighting is often called factual explanations.

37 Inspired by human scaffolding behavior, our main contribution is extending the common preference-  
38 based interface to allow the human to explain his decision more. More specifically, we allow the  
39 human to select timesteps of both trajectories (of the pair) to annotate which steps they deem important  
40 for his decision. This can be seen as a factual explanation. Our method then uses these explanations  
41 as additional training input. We archive this by generating an explanation of the reward model  
42 through gradients (e.g., saliency-based explanations) and comparing the generated explanations with  
43 the human explanations. This additional loss can be transparently integrated into the reward model  
44 learning of any preference-based reinforcement learning framework.

## 45 2 Related Work

46 Multiple works have attempted to extend the pipeline of Human-in-the-Loop RL to include more  
47 than (evaluative) feedback. Guan et al. [2020] implemented an HRL method that allows the user to  
48 highlight important areas of an image in addition to evaluative feedback. Our implementation can  
49 be used with any preference-based method (which can be used with any RL algorithm), while they  
50 designed a custom HRL method. Additionally, our proposed solution has the advantage that it is not  
51 restricted to an image-based state space but can deal with any state space.

52 Mahmud et al. [2023] showed that human explanation on individual feature space level can help  
53 to create better rewards models from an offline dataset of demonstrations. Again, their method is  
54 limited to cases where humans can understand the feature space. And secondly, their work is on  
55 offline learning, while we learn online.

56 Basu et al. [2018] allow the human to select which features influence their decision between two  
57 trajectories. While they show that their method can lead to increased learning performance, it is  
58 limited to a linear reward model, which limits its application quite drastically. Additionally, their  
59 work is limited in that they rely on complete trajectories (instead of segments in the case of PBRL).  
60 From the human viewpoint, the work of Basu et al. [2018] is the closest comparison to ours.

61 Saran et al. [2021] improved the efficiency of Imitation Learning by collecting human gaze data,  
62 which is then used in an auxiliary loss constructed from the network attention and the human. Their  
63 idea of including additional information through an auxiliary loss, which is constructed through  
64 gradients of the network, is similar to the approach. But their focus is more on state regions, not  
65 timesteps, and the different settings of imitation learning.

66 Karalus and Lindner [2022] used counterfactuals (instead of factuals like our method) to increase the  
67 learning speed of TAMER (Knox and Stone [2008] alternative to preference-based RL for evaluative  
68 feedback). They reported increased learning performance, but their methods are for different types of  
69 human explanations and formulated in a different framework.

70 Gajcin et al. [2023] lets users select part of a trajectory that they do not like. With this information,  
71 they perform reward shaping. While they allow humans to select parts of trajectories that they don't  
72 like, their annotations focus on either states or actions, not timesteps. Additionally, their setting of  
73 reward shaping still assumes the existence of an environmental reward and does not allow the training  
74 of agents purely from human feedback.

75 Wu et al. [2024] let the user select text parts (i.e., parts of their trajectory) to specify if that part is  
76 relevant for one of their pre-defined classes (Relevance, Factuality, etc..) While their work is solely  
77 focused on LLM/RLHF, the concept of selecting parts of the input trajectory is similar to our work.  
78 But this additional information is then used to train a pre-selected hard-coded set of reward models,  
79 which makes the transfer to other domains different. While our method is much more flexible and  
80 requires less oversight.

81 Therefore, the opportunity to extend the setting to allow humans to give additional explanations has  
82 been investigated in different settings but not in the preference-based framework. Additionally, most  
83 of the work requires hand-crafted algorithms, while our contribution is much more independent of the  
84 concrete PBRL method. However, most methods reported increased learning speeds when including  
85 additional human annotations (in any form); therefore, it's realistic to expect changes in learning  
86 speed from our work.

87 Metz et al. [2023] proposed an interface with multitudes of different annotation options for collecting  
88 human preferences over trajectories. Their focus was purely on the interface design, which is not  
89 connected to any real algorithm. To our knowledge, an explanatory annotation of timesteps inside  
90 trajectories is impossible in their theoretical interface because they drew heavy inspiration from  
91 existing research. This shows the importance of showcasing the technical feasibility of possible

92 feedback interface avenues so they can be included in future human-computer-interaction research.  
 93 Casper et al. [2023] specified "fine-grained feedback" as a possibility for problems where precise  
 94 information is necessary to solve the task.

95 In settings other than PBRL, the concept of using gradient-based explanations to enhance training  
 96 has been applied successfully. Saliency-based methods have been used in other contexts to speed  
 97 up the learning process or force the learner to learn the right reasons. Common examples where  
 98 gradient-based explanations have been used (in supervised learning) to accelerate learning are  
 99 ExpectedGradients (Erion et al. [2021]), Right-for-the-right-reasons (Ross et al. [2017] or Ismail et al.  
 100 [2021]). These ideas are not yet been explored in PBRL.

### 101 3 Background

102 We consider a setting wherein the sequential interactions between an agent and its environment are  
 103 formally characterized by a Markov Decision Process (MDP), a mathematical framework comprising  
 104 states, actions, and transitions. Within this formalism, the system evolves discretely at time intervals  
 105  $t$ . Within the episodic framework, the agent’s involvement persists until reaching a terminal time step  
 106  $T$ , signifying the conclusion of an episode. This temporal progression is succinctly represented by the  
 107 trajectory, expressed as the ordered sequence  $(s_1, a_1), \dots, (s_T, a_T)$ , encapsulating the observation-  
 108 action pairs throughout an episode.

109 In RL the agent experiences rewards at each time step. In our Human-in-the-Loop context, we abstain  
 110 from presuming access to such rewards. Instead, we introduce a human overseer controlling the  
 111 agent’s task intention through two distinct feedback channels. First, their preferences (as generally in  
 112 PBRL) and second, their explanations (our contribution).

#### 113 3.1 Preference-based Reinforcement Learning

114 The main goal of preference-based reinforcement learning (PBRL, Christiano et al. [2017]) is to learn  
 115 a reward function, denoted as  $\hat{r}_\psi$ , from a set of expressed segment pairs. These pairs of segments  
 116 are collected from the human throughout the training process. Within this framework, a segment,  
 117  $\sigma$ , is defined as a sequence of states and actions,  $\{s_k, a_k, \dots, s_{k+H}, a_{k+H}\}$ . Preferences, denoted as  
 118  $y$ , are elicited for segments  $\sigma_0$  and  $\sigma_1$ , with  $y$  representing a distribution indicating the preferred  
 119 segment, i.e.,  $y \in \{(0, 1), (1, 0), (0.5, 0.5)\}$ . This evaluative judgment is recorded in a dataset  $D$  as a  
 120 triple  $(\sigma_0, \sigma_1, y)$ . This dataset of binary preference can be then used to train a reward model with the  
 121 Bradley-Terry model (Bradley and Terry [1952]):

$$P_\psi[\sigma_1 \succ \sigma_0] = \frac{\exp \sum_t \hat{r}_\psi(s_{1t}, a_{1t})}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}_\psi(s_{it}, a_{it})} \quad (1)$$

122 Here,  $\sigma_i \succ \sigma_j$  denotes that segment  $i$  is preferred to segment  $j$ . This formulation states that the  
 123 probability of preferring a segment exponentially depends on the reward function’s sum over the  
 124 segment. While  $\hat{r}_\psi$  itself is not inherently a binary classifier, the learning procedure is akin to binary  
 125 classification, where a supervisor supplies labels  $y$ . To learn the reward function, instantiated as a  
 126 neural network with parameters  $\psi$ , the following loss is minimized:

$$\mathcal{L}_{\text{preference}} = -\mathbb{E}_{(\sigma_0, \sigma_1, y) \sim \mathcal{D}} [y_0 \log P_\psi[\sigma_0 \succ \sigma_1] + y_1 \log P_\psi[\sigma_1 \succ \sigma_0]] \quad (2)$$

127 Even in PBRL the final goal is to learn a policy  $\pi$  that maximizes the expected cumulative reward.  
 128 The reward function  $\hat{r}_\psi$  obtained from preferences is a surrogate to the true, unknown reward function.  
 129 Therefore  $\hat{r}_\psi$  is utilized instead of the actual reward signal while optimizing the policy. The policy  
 130 can be learned through standard RL algorithms like PPO (Schulman et al. [2017]) or SAC (Haarnoja  
 131 et al. [2018]).

#### 132 3.2 Training Guided by Saliency-Based Methods

133 Saliency-based explanation methods allow an understanding of neural network decision-making  
 134 processes, offering a quantitative means to assess the impact of input features on model predictions.  
 135 Consider a neural network  $f_\theta$  with parameters  $\theta$ , trained on a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ , where

136  $x_i$  denotes input samples and  $y_i$  represents corresponding labels. To calculate a saliency-based  
 137 explanation (or saliency map) for an input  $x_i$ , a normal loss between the prediction of the network  
 138  $f_\theta$  and the label  $y_i$  is calculated. That loss then allows us to calculate the gradients of the input  $x_i$   
 139 (with respect to the before-mentioned loss). These gradients are then used as an explanation of the  
 140 prediction.

141 One notable saliency-based method is SmoothGrad (Smilkov et al. [2017]), designed to stabilize  
 142 saliency maps. The smoothed saliency map  $S_{\text{SmoothGrad}}$  is computed as the average gradient over  
 143  $N_{\text{smooth}}$  perturbed samples:

$$S_{\text{SmoothGrad}} = \frac{1}{N_{\text{smooth}}} \sum_{i=1}^{N_{\text{smooth}}} \nabla_x f_\theta(x + \epsilon_i) \quad (3)$$

144 Here,  $\epsilon_i$  represents the perturbation, and  $\nabla_x f_\theta(x + \epsilon_i)$  is the gradient of the model’s prediction with  
 145 respect to the perturbed input. This approach introduces controlled noise to the input, providing a more  
 146 stable saliency map. The resulting map highlights the regions crucial for the network’s predictions,  
 147 contributing to a nuanced understanding of feature importance in neural network decision-making.  
 148 Similar methods like IntegratedGradients Sundararajan et al. [2017] calculate multiple explanations  
 149 with respect to a baseline and average the results to increase stability.

150 Saliency-based training methods (Ismail et al. [2021]), such as ExpectedGradients (Erion et al. [2021])  
 151 and Right-for-the-Right-Reasons (Ross et al. [2017]), allow the shaping of the training dynamics  
 152 of neural networks. These approaches utilize saliency information to refine the training process,  
 153 enhancing model interpretability and predictive performance. Right-for-the-Right-Reasons emphasize  
 154 salient features contributing to correct predictions during training, steering the model towards more  
 155 meaningful representations. This is done by adding additional terms to the loss term. Instead of  
 156 a single loss term between the prediction’s distance from the target (i.e., being right), a second  
 157 (weighted) term is added to constrain the distance between a saliency-based explanation and a given  
 158 explanation (i.e., the right reason).

159 ExpectedGradients Erion et al. [2021] allows to guide the training even in cases where no ground-truth  
 160 explanations are available. This is done by imposing priors on the structure of the explanation. For  
 161 example, explanations should be either locally-smooth or sparse regarding features. These saliency-  
 162 guided training strategies embody a synergy between interpretability and maintaining predictive  
 163 accuracy.

## 164 4 Implementation

165 Our main contribution is the extension of the feedback possibilities in PBRL for humans. In addition to  
 166 their preferences between the two trajectories, the human can now explain which timesteps they deem  
 167 essential for their decision. This evaluative judgment, accompanied by detailed binary explanations  
 168 provided by the human over timesteps, is recorded the dataset  $\mathcal{D}$  as a quintuple  $(\sigma_0, \sigma_1, y, e_1, e_2)$ .  
 169 Here,  $e_1$  and  $e_2$  directly represent binary vectors of human-provided explanations, each having the  
 170 same length as the corresponding segments  $\sigma_0$  and  $\sigma_1$ . These binary vectors,  $e_1$  and  $e_2$ , serve to  
 171 elucidate the significance of individual timesteps within their respective segments, providing binary  
 172 indications (0 or 1) of the human’s perceived importance at each step. Figure 1 shows an overview of  
 173 our approach.

### 174 4.1 Reward Learning with annotated preferences

175 Incorporating preference-based learning alongside annotations necessitates introducing two additional  
 176 steps within the learning process. First, we must predict an explanation  $\hat{e}$  and then compare this  
 177 generated explanation with the human explanation  $e$ .

178 Initially, an explanation of the existing prediction (of the reward model) concerning the current  
 179 trajectory pair is formulated. The generation of this explanation leverages established saliency-based  
 180 methods, and we employ the SmoothGrad technique (see Equation 3) in our implementation. Here,  
 181 a set of  $n_{\text{smooth}}$  perturbations is generated for a given input trajectory by sampling from a normal  
 182 distribution with mean and standard deviation parameters derived from the input trajectory. On these  
 183 perturbations, we then calculate the saliency-based explanations.

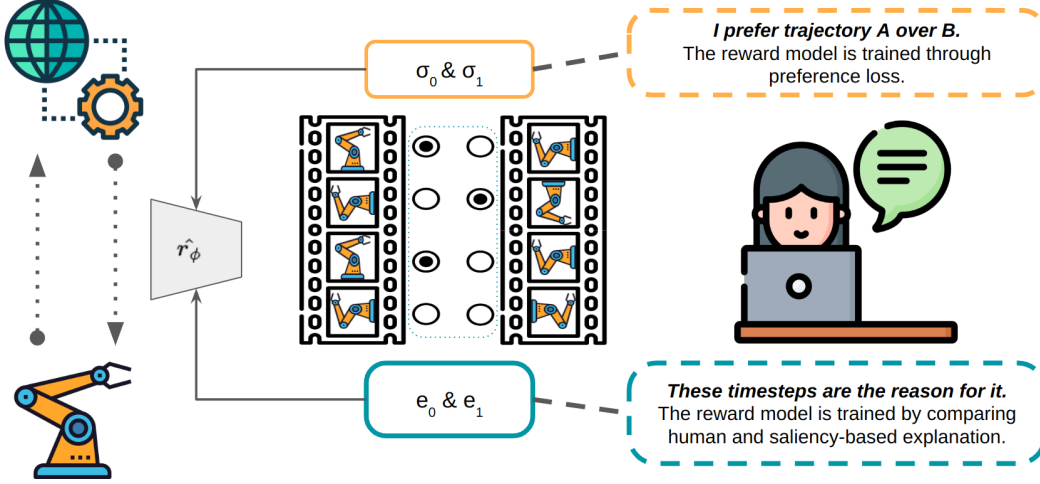


Figure 1: Overview of our approach. Left: The agent optimizes his policy with respect to the trained proxy reward model. Right: The proxy reward model is trained from human preferences over trajectories and the importance of each timestep in the trajectory. Humans provide both.

$$\xi(\sigma) = \left[ \sum \left| \frac{1}{N_{\text{smooth}}} \sum_i^{N_{\text{smooth}}} \nabla_{s,a} \hat{r}_\psi(s + \epsilon_i, a + \epsilon_i) \right| \right] \forall s, a \in \sigma \quad (4)$$

184  $\epsilon$  represents the perturbation of the input. We then predict the outcome (reward) with the reward  
 185 model  $\hat{r}_\psi$  and calculate the gradient with respect to the input, i.e., the state and action. This done for  
 186 all timesteps  $t$  in the segment  $\sigma$ .

187 Since the human explanations ( $e_1$  and  $e_2$ ) are binary vectors with the length of the segments ( $\sigma_0, \sigma_1$ ),  
 188 but the saliency-based explanations  $\hat{e}$  are of shape  $(s, a)_t$  (a concatenation of state and actions space  
 189 for each timestep in the segment). Therefore, the generated explanations must be transformed from  
 190 the state & action frame space to a step-level explanation. To achieve this, we sum the absolute  
 191 values over each timestep to transform the frame-based values into a step-level explanation. These  
 192 unnormalized values then serve as logits for the binary-cross entropy loss, aligning the calculated  
 193 explanation with the provided annotation, the latter specified by the user’s annotations.

194 The annotation loss  $L_{\text{annotation}}$  is the standard binary multi-class cross-entropy between our explanation-  
 195 logits and the true explanation:

$$L_{\text{annotation}} = \frac{1}{2N} \sum_{e, \sigma \in \{(e_0, \sigma_0)(e_1, \sigma_1)\} \in \mathcal{D}} [-e \cdot \log(f(\xi(\sigma))) + (1 - e \cdot \log(1 - f(\xi(\sigma))))] \quad (5)$$

196 where  $e$  is the explanation of the user (for the human preference  $y$ ) and  $\xi(\sigma)$  generates the explanation  
 197 on a trajectory step level.  $N$  is the total number of preferences (i.e., a pair of segments) in the database  
 198  $\mathcal{D}$ . The user’s preferences are  $e$ , and our generated preferences are transformed with  $\xi(\sigma)$ . Since  $\xi$   
 199 produces normalized logits, we apply the sigmoid function for  $f$ . To relax the learning constraints,  
 200 we employ label smoothing for the human explanations  $e$  to relax the constraints of having the correct  
 201 explanations.

202 Theoretically, we could allow the user to give specific annotations, i.e., highlighting single entries  
 203 in the state. Still, there are reasons to limit the annotations’ granularity to the trajectory’s timestep  
 204 level. First, the timestep level of the trajectory is quite application-independent and does not require  
 205 the human to fully understand lower levels (like the state and action space). For example, in many  
 206 robots, the state space includes a mix of different sensors (LIDAR, Camera, etc.), which can be hard  
 207 to understand fully for humans. Lower levels also drastically increase the mental load on humans due  
 208 to the massive number of options. Therefore, we want to show in this work that even explanations on  
 209 the timestep level contain enough information to increase learning speed.

## 210 4.2 Structural Loss

211 Like ExpectedGradients, we want to constrain the structure of our generated explanation and impose  
212 a structural loss on the derived explanation. Since most human explanations are sparse, we apply  
213 this prior to our generated explanation and force the generated explanations to reflect this sparse  
214 characteristic. This structural loss is realized by applying an L1 norm to the calculated explanation.

215 Mathematically, the structural loss  $L_{\text{structural}}$  is formulated as:

$$L_{\text{structural}} = \frac{1}{2N} \sum_{\sigma \in \{\sigma_0, \sigma_1\} \in \mathcal{D}} \|\xi(\sigma)\|_1 \quad (6)$$

216 where  $N$  denotes the total number of trajectories in  $\mathcal{D}$ , and  $\xi$  generates a step level explanation for  
217 the segment  $\sigma$ . Incorporating this structural loss contributes to the regularization of the explanation,  
218 fostering a more coherent and interpretable representation.

219 To derive our reward model  $\hat{r}_\phi$ , we integrate all three distinctive loss terms: the primary preference  
220 loss (Equation 2), the annotation loss (Equation 5), and the structural loss (Equation 6), merging them  
221 into a unified loss term. The incorporation of the additional losses is accompanied by weightings ( $\alpha_1$   
222 and  $\alpha_2$ ):

$$L_{\text{total}} = L_{\text{preference}} + \alpha_1 \cdot L_{\text{annotation}} + \alpha_2 \cdot L_{\text{structural}} \quad (7)$$

223 This loss formulation ensures the simultaneous consideration of preference-based learning, explana-  
224 tion fidelity, and structural regularization in the training of the final reward model  $\hat{r}_\phi$ .

## 225 4.3 Policy Learning

226 To learn an actual policy from the learned reward function, we use the same setup as PEBBLE Lee  
227 et al. [2021a], in which the offline learning algorithm SAC is leveraged. Like in PEBBLE, we update  
228 the replay buffer every time the reward function is updated. This allows the agent to always learn  
229 from the latest experience. Since we extend from PEBBLE, it’s also the choice of our baseline for the  
230 evaluation.

## 231 5 Evaluation

232 We evaluate our changes in three common robotics environments, Walker, HalfCheetah, and Hopper  
233 of the mujoco suite. Most of the hyperparameters are taken from our baseline PEBBLE (Lee et al.  
234 [2021a]), which we then briefly tuned only on the baseline (since we argue that our work should be  
235 viewed as an extension to the existing method, we believe it should not deserve a full hyperparameter  
236 search). The ballpark for the hyperparameter of our extension (i.e., the weight of additional loss  
237 terms) was also taken from similar methods (Ross et al. [2017]) and briefly sanity-checked. The  
238 agent trains for 2 million steps in each environment, with a total feedback budget of 700 comparisons.  
239 A full overview of all hyperparameters can be found in Appendix A

240 While we are learning in an HRL setting, where we don’t have access to the true reward function,  
241 we still use the true reward function for the evaluation. To provide a solid measurement of the  
242 learning progress, we perform every 10k environment steps performance measurement of the agent.  
243 This measurement consists of 5 episodes in a newly seeded (with a different random seed than the  
244 current iteration) environment. At each measurement, we collect total true environmental rewards  
245 and calculate the mean between the 5 measurement episodes. With this setting, we collect 200  
246 measurements throughout the whole training process (of a single evaluation run). This extensive  
247 collection of measurements gives us a more stable estimation of the true performance (of the agent)  
248 throughout training than the normally used rollout mean rewards. These collected measurements  
249 are then aggregated and normalized between zero and one. We use the same procedure proposed by  
250 Agarwal et al. [2021] to calculate mean scores and confidence intervals via bootstrapping over 10  
251 runs in each condition.

252 We opted to create a synthetic oracle instead of actual humans to allow for an evaluation with enough  
253 statistical significance. For the comparison, we used the environmental ground truth rewards (since in

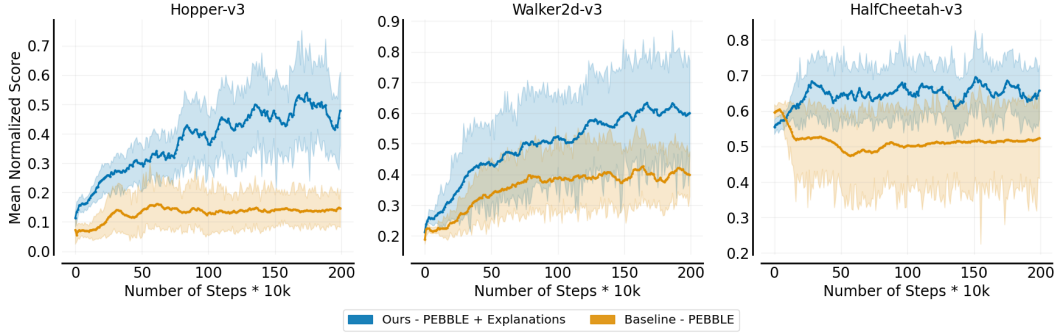


Figure 2: The mean throughout the training process in each environment. Shaded areas represent the confidence interval. Higher is better.

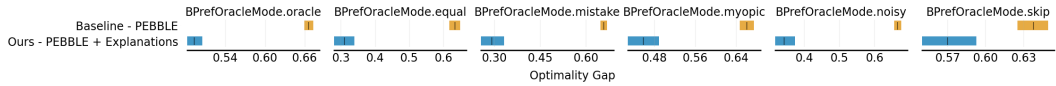


Figure 3: Optimality Gap (how fast a run converges to the ideal score, lower is better) under different human irrationality. The vertical line visualizes the mean optimality gap, and the bars represent the confidence interval.

254 DM-control, it exists here); for the trajectory step-level explanation, we resorted to an XAI technique  
 255 called IntegratedGradients (Sundararajan et al. [2017]) to extract the explanations from a pre-trained  
 256 policy.

257 The random seed is separated between the environment dynamics and the network/agent dynamics.  
 258 We select the same starting seeds for the comparisons between our work and the baseline (to reduce  
 259 the impact of "favorable" seeds).

260 Since the training with a perfect oracle is not a realistic setting, we test the stability of our extension  
 261 properly and use the five irrationalities of human teachers as defined by the standard B-Pref benchmark  
 262 (see Lee et al. [2021b]). These different common failure patterns each represent different "sub-  
 263 optimality" that can occur in human feedback. We also train with the perfect case (the oracle case).  
 264

## 265 5.1 Results

266 As visible in Figure 2, we can achieve better mean rewards (the mean of 10 runs, each point  
 267 is measured in 5 evaluation episodes, as described in section 5) results than the baseline with  
 268 our extension in all environments with the perfect oracle. While in the Walker and HalfCheetah  
 269 environments, the mean (bold line) is higher, the confidence interval (shaded areas) still overlaps;  
 270 therefore, we can't extensively say that our methods are conclusively better. Only in the Hopper case  
 271 can we conclusively answer that the human explanation helps the reward model learn faster. While  
 272 not pictured, both methods converge to the same reward in the long run. However, our method's  
 273 major motivation and driving force was increasing learning speed, especially in the earlier phases  
 274 when the human is still in the loop.

275 Since we want to highlight the increase in learning speed, we use the optimality gap as our primary  
 276 metric for further evaluation since it condenses the learning speed into a single, more understandable  
 277 metric. In Figure 3, we use the optimality gap to compare our modifications in the different irrational-  
 278 ities defined by BPref. Notable is the significance (indicated by no overlapping bars, as defined in  
 279 Agarwal et al. [2021]) between our work and the baseline in quite a few cases. These results confirm  
 280 that our extension leads to better learning performance, even in settings where the human feedback  
 281 follows a more realistic pattern.

## 282 6 Discussion

283 **Additional mental load** Labelling trajectories not only with the binary preference but also with the  
284 annotation of which parts are better/worse does, of course, require more mental effort than pure  
285 preference-based (or other evaluative approaches). In contrast, we showed that these annotations  
286 speed up the learning process, and therefore, less total human interaction is needed. We believe this  
287 trade-off has to be evaluative on a case-by-case basis since it's dependent on multiple factors like the  
288 environment, the learning algorithm, the feedback interface, and the expertise and type of trainers.  
289 Additionally, we want to highlight that our annotations might relieve some of the frustrations human  
290 trainers often have when using a very restricted feedback interface (binary in the case of preference-  
291 based learning). Our methods give the human trainer more avenues to express his feedback, which  
292 should increase user experience.

293 **Granularity of the annotations** While it's technically feasible to annotate whole states/actions in  
294 a trajectory and annotate specific parts of the state space directly, we refrain from evaluating this  
295 scenario. We believe preference-based learning shines in settings where the human can't highlight  
296 single points in the state-space because of the complexity of the state space. We believe in areas  
297 where its possible for human to understand single points in the state-space, methods like TAMER  
298 Knox and Stone [2008] are probably a better alternative.

299 **Evaluation with a simulated oracle** Throughout this work, we evaluated our modification with  
300 a synthetic oracle instead of an actual human. While this is a weakness of the evaluation, it also  
301 allowed us to collect more runs (with different random seeds) and perform different ablations. This  
302 increase makes the statistical analysis more sound. We believe this benefits the research community  
303 more than a human study with a few participants, which would not allow us to show our gains with  
304 (statistical) confidence.

## 305 Conclusion

306 With this work, we showed how preference-based reinforcement learning could be extended so  
307 humans can give not only preferences but also explanations for their decisions. In experiments, we  
308 showed that with these explanations, we can increase the performance of current state-of-the-art PBRL  
309 even further. Not only in the ideal case but also under a variety of different human irrationalities, our  
310 extension either increases the performance or does not reduce it.

311 Our method, as presented here, is meant as a first demonstration, favoring simplicity over sophis-  
312 tication. A significant positive attribute of PBRL is that its RL part is mainly decoupled from the  
313 human-in-the-loop, which means that future advantages in RL, which bring new algorithms, can be  
314 easily included. Our extensions continue that attribute since we only require changes to the reward  
315 model, not the agent.

316 Overall, we showed potential in opening up the interface in human-in-the-loop/preference-based  
317 RL to include more natural and human-like feedback mechanisms instead of relying only on binary  
318 feedback.

## 319 References

- 320 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare.  
321 Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Information*  
322 *Processing Systems*, 2021.
- 323 Saleema Amershi, Maya Cakmak, William Bradley Knox, and Todd Kulesza. Power to the people:  
324 The role of humans in interactive machine learning. *AI Magazine*, 35(4):105–120, Dec. 2014.
- 325 Chandrayee Basu, Mukesh Singhal, and Anca D. Dragan. Learning from richer human guidance:  
326 Augmenting comparison-based learning with feature queries. *2018 13th ACM/IEEE International*  
327 *Conference on Human-Robot Interaction (HRI)*, 2018.
- 328 Serena Booth, W. Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi. The  
329 perils of trial-and-error reward design: Misdesign through overfitting and invalid task specifications.  
330 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.



- 331 Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method  
332 of paired comparisons. *Biometrika*, 1952. ISSN 00063444.
- 333 Stephen Casper, Xander Davies, and et al. Open problems and fundamental limitations of reinforce-  
334 ment learning from human feedback. *Transactions on Machine Learning Research*, 2023.
- 335 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
336 reinforcement learning from human preferences. *Advances in neural information processing*  
337 *systems*, 30, 2017.
- 338 Gabriel G. Erion, Joseph D. Janizek, Pascal Sturmfels, Scott M. Lundberg, and Su-In Lee. Learning  
339 explainable models using attribution priors. *Nature machine intelligence*, 3, 2021.
- 340 Jasmina Gajcin, James McCarthy, Rahul Nair, Radu Marinescu, Elizabeth M. Daly, and Ivana  
341 Dusparic. Iterative reward shaping using human feedback for correcting reward misspecification.  
342 *Proceedings of the 2023 European Conference on Artificial Intelligence*, abs/2308.15969, 2023.
- 343 L. Guan, Mudit Verma, Sihang Guo, Ruohan Zhang, and Subbarao Kambhampati. Widening  
344 the pipeline in human-guided reinforcement learning with explanation and context-aware data  
345 augmentation. In *Neural Information Processing Systems*, 2020.
- 346 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
347 maximum entropy deep reinforcement learning with a stochastic actor. *International conference*  
348 *on machine learning*, pages 1861–1870, 2018.
- 349 Aya Abdelsalam Ismail, Hector Corrada Bravo, and Soheil Feizi. Improving deep learning inter-  
350 pretability by saliency guided training. *Advances in Neural Information Processing Systems*, 34,  
351 2021.
- 352 Jakob Karalus and Felix Lindner. Accelerating the learning of TAMER with counterfactual expla-  
353 nations. In *2022 IEEE International Conference on Development and Learning (ICDL)*, pages  
354 362–368. IEEE, 2022.
- 355 W. Bradley Knox and Peter Stone. TAMER: Training an Agent Manually via Evaluative Reinforce-  
356 ment. In *IEEE 7th International Conference on Development and Learning*, August 2008.
- 357 Kimin Lee, Laura Smith, and Pieter Abbeel. PEBBLE: Feedback-efficient interactive reinforcement  
358 learning via relabeling experience and unsupervised pre-training. In *International Conference on*  
359 *Machine Learning*, 2021a.
- 360 Kimin Lee, Laura Smith, Anca Dragan, and Pieter Abbeel. B-Pref: Benchmarking preference-based  
361 reinforcement learning. In *Thirty-fifth Conference on Neural Information Processing Systems*  
362 *Datasets and Benchmarks Track (Round 1)*, 2021b.
- 363 Saaduddin Mahmud, Sandhya Saisubramanian, and Shlomo Zilberstein. Explanation-guided reward  
364 alignment. In *Proceedings of the Thirty-Second International Joint Conference on Artificial*  
365 *Intelligence, IJCAI-23*. International Joint Conferences on Artificial Intelligence Organization,  
366 2023.
- 367 Yannick Metz, David Lindner, Raphael Baur, Daniel A. Keim, and Mennatallah El-Assady. RLHF-  
368 Blender: A configurable interactive interface for learning from diverse human feedback. *Interactive*  
369 *Learning with Implicit Human Feedback Workshop at ICML 2023*, abs/2308.04332, 2023.
- 370 Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons:  
371 Training differentiable models by constraining their explanations. *Proceedings of the Twenty-Sixth*  
372 *International Joint Conference on Artificial Intelligence, IJCAI-17*, 2017.
- 373 Akanksha Saran, Ruohan Zhang, Elaine Schaertl Short, and Scott Niekum. Efficiently guiding  
374 imitation learning algorithms with human gaze. *International Conference on Autonomous Agents*  
375 *and Multiagent Systems (AAMAS)*, abs/2002.12500, 2021.
- 376 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
377 optimization algorithms. *ArXiv*, abs/1707.06347, 2017.

- 378 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad:  
379 removing noise by adding noise. *Workshop on Visualization for Deep Learning, ICML, 2017*.
- 380 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
381 *International conference on machine learning*. PMLR, 2017.
- 382 Janneke van de Pol, Monique Volman, and Jos Beishuizen. Scaffolding in teacher–student interaction:  
383 A decade of research. *Educational Psychology Review*, 22, 09 2010.
- 384 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith,  
385 Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for  
386 language model training. *Advances in Neural Information Processing Systems*, 36, 2024.

## 387 **A Hyperparameters**

388 We use the following hyperparameters for our experiments. Most of the common hyperparameters  
389 are from Lee et al. [2021a] or Christiano et al. [2017].

### 390 **A.1 Compute Power**

391 Most of the experiments were performed on a single NVIDIA A100 with a single seed run in  
392 environments that took roughly 4 hours. Many runs were calculated in parallel on the same machine  
393 to increase speed.

<b>Hyperparameter</b>	<b>Value</b>
<b>SAC</b>	
Learning rate	0.0003
Batch size	512
Buffer size	1000000
$\tau$	0.005
$\gamma$	0.99
<b>Reward Model</b>	
Ensemble size	3
Number of hidden units	300
Number of hidden layers	3
Learning rate	5e-4
Batch size	32
Adam $\beta'$ s	(0.9, 0.9)
AdamW weight decay	0.05
<b>Our</b>	
$\alpha_1$	0.25
$\alpha_2$	0.1
Smoothgrad number of perturbations	16
Smoothgrad noise	0.01
<b>Feedback</b>	
Maximum number of feedback	700
Feedback per interval	70
Feedback is given every X steps	20000
Maximum segment size	50
Segment sampling scheme	Disagreement

Table 1: Hyperparameters