

LEVERAGING SUB-OPTIMAL DATA FOR HUMAN-IN-THE-LOOP REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

To create useful reinforcement learning (RL) agents, step zero is to design a suitable reward function that captures the nuances of the task. However, reward engineering can be a difficult and time-consuming process. Instead, human-in-the-loop (HitL) RL methods hold the promise of learning reward functions from human feedback. Despite recent successes, many of the HitL RL methods still require numerous human interactions to learn successful reward functions. To improve the feedback efficiency of HitL RL methods (i.e., require less human interaction), this paper introduces Sub-optimal Data Pre-training, SDP, an approach that leverages reward-free, sub-optimal data to improve scalar- and preference-based HitL RL algorithms. In SDP, we start by pseudo-labeling all low-quality data with the minimum environment reward. Through this process, we obtain reward labels to pre-train our reward model *without* requiring human labeling or preferences. This pre-training phase provides the reward model a head start in learning, enabling it to recognize that low-quality transitions should be assigned low rewards. Extensive experiments with both simulated and human teachers reveal that SDP can at least meet, but often significantly improve, state-of-the-art HitL RL performance across a variety of simulated robotic tasks.

1 INTRODUCTION

In reinforcement learning (RL), an agent’s objective is to interact with an environment and maximize its total (discounted) expected reward. The reward hypothesis further maintains that a well-specified reward function is sufficient for an agent to learn to solve a task (Sutton & Barto, 2018). However, defining a reward function that precisely captures all task complexities is often tedious and non-trivial (Booth et al., 2023). There have been notable examples of reward misspecification, in which RL agents discovered and exploited unintended shortcuts in the reward function (Skalse et al., 2022). One notorious example is the CoastRunners game, in which the goal should be to finish a boat race as fast as possible — an RL agent instead gained the most reward by spinning its boat in a circle despite concurrently catching on fire and crashing into other boats (Clark & Amodei, 2016).

A promising alternative is to learn reward functions directly from human feedback. In this paradigm, humans can provide feedback in the form of preferences or scalar signals, which can then be used to learn a reward function that is consistent with human desires (Daniel et al., 2014; Christiano et al., 2017). Despite recent progress, existing preference- and scalar-based RL methods still suffer from high human labeling costs that can require thousands of human queries to learn an adequate reward function (Christiano et al., 2017). Prior work attempts to mitigate this issue through several mechanisms, including active learning (Lee et al., 2021a), data augmentation (Park et al., 2022), semi-supervised learning (Park et al., 2022), and meta-learning (Hejna III & Sadigh, 2023).

Alternatively, our work draws on recent advances in offline RL that have demonstrated the value of low-quality data (Yu et al., 2021). However, its potential in HitL RL remains unexplored. As low-quality data is often readily accessible or easy to obtain, this work addresses this gap by asking the question:

Can we leverage sub-optimal, unlabeled data to improve learning in HitL RL methods?

To that end, we present Sub-optimal Data Pre-training, *SDP*, a tool for HitL RL algorithms to increase human feedback efficiency. SDP leverages sub-optimal trajectories by pseudo-labeling all

054 transitions with the minimum environment reward. The now pseudo-labeled sub-optimal data serves
 055 two purposes. First, we pre-train a regression-based reward model by applying standard supervised
 056 learning to minimize the mean squared loss. Intuitively, this pre-training provides the reward model
 057 a head start, biasing it towards assigning lower reward values to these low-quality transitions. Sec-
 058 ond, we initialize the RL agent’s replay buffer with the sub-optimal data and make learning updates
 059 to the RL agent. This process changes the RL agent’s policy and provides different behaviors for the
 060 human to provide feedback on (relative to learning with no initial sub-optimal data). This ensures
 061 that when the human teacher provides feedback, their time is used efficiently, avoiding redundant
 062 feedback on the existing sub-optimal data. Afterward, we follow the standard preference- or scalar-
 063 based RL protocol.

064 This paper’s core contribution is showing that we can harness the availability of low-quality, reward-
 065 free data for HitL RL approaches by pseudo-labeling it with minimum rewards and treating it as a
 066 prior for learning reward models. We first validate the utility of SDP in extensive simulated teacher
 067 experiments, combining it with four scalar- and preference-based RL algorithms. These experiments
 068 show that SDP significantly improves the efficiency of feedback in complex tasks from both the
 069 DeepMind Control (Tassa et al., 2018) and Meta-World (Yu et al., 2020) suites. Crucially, we further
 070 highlight the real-world applicability of SDP by demonstrating its success with human teachers in a
 071 16-person user study. Overall, this work takes an important step toward considering how HitL RL
 072 approaches can take advantage of readily-available sub-optimal data.

073 2 RELATED WORK

074 **Human-in-the-Loop RL** Several approaches in HitL RL allow agents to leverage human feedback
 075 to adapt or learn new behavior. Learning from demonstration is one such methodology that allows
 076 a human to provide examples of desired agent behavior (Argall et al., 2009). Human demonstration
 077 data has been used to shape the environment’s reward function (Brys et al., 2015), develop a reward
 078 function from scratch (Abbeel & Ng, 2004), and bias the agent’s policy towards certain actions
 079 (Taylor et al., 2011). Although demonstrations can be a rich source of feedback, they are often
 080 expensive to obtain and may require domain experts (Dragan & Srinivasa, 2012).
 081

082 Another approach is learning from preference-based feedback where a teacher provides preferences
 083 between two or more sets of agent behavior (Christiano et al., 2017). Preference learning has been
 084 popularized in recent years as it can require less effort and expertise compared to providing demon-
 085 strations. To further reduce the amount of human interaction required, several strategies have been
 086 introduced. This has included combining preferences with demonstrations (Ibarz et al., 2018; Bıyık
 087 et al., 2022), unsupervised pre-training (Lee et al., 2021a), bi-level optimization (Liu et al., 2022),
 088 semi-supervised learning (Park et al., 2022), data augmentation (Park et al., 2022), uncertainty-based
 089 exploration (Liang et al., 2022), meta-learning (Hejna III & Sadigh, 2023), and active learning ap-
 090 proaches (Hu et al., 2024). Despite its popularity, some argue that comparison feedback might not
 091 capture the full intricacies of human preferences, as oftentimes the human is limited to choosing
 092 between two options (Daniel et al., 2014; White et al., 2024).
 093

094 As a result, another body of work focuses on learning from scalar feedback where human teachers
 095 can provide scalar signals to evaluate an agent’s behavior (Knox & Stone, 2009; Griffith et al., 2013;
 096 Loftin et al., 2016; MacGlashan et al., 2017; White et al., 2024). Several works use scalar feedback
 097 to either learn a reward model (Daniel et al., 2014; Cabi et al., 2020) or an action-value function
 098 (Knox & Stone, 2009; 2013; Warnell et al., 2018) via regression.

099 **Learning from Sub-Optimal Data** SDP aims to leverage sub-optimal data for scalar- and
 100 preference-based RL algorithms. However, learning from low-quality data or negative examples has
 101 been applied in other areas of RL and imitation learning (Chen et al., 2021; Tangkaratt et al., 2021).
 102 In standard RL, several works use sub-optimal demonstrations to initialize a policy (Taylor et al.,
 103 2011; Hester et al., 2018; Gao et al., 2019). In goal-conditioned RL, Hindsight-Experience-Replay
 104 uses failed episodes by treating them as a success with respect to a different goal (Andrychowicz
 105 et al., 2017). In inverse reinforcement learning (IRL), Shiarlis et al. (2016) proposed a constrained
 106 optimization formulation that can accommodate both successful and failed demonstrations. Brown
 107 et al. (2019) makes use of ranked demonstrations to learn a reward function in IRL. Later work in
 IRL automatically generates ranked trajectories by adding increasing amounts of noise to a learned

policy (Brown et al., 2020). Lastly, in offline RL, Singh et al. (2020) leverages sub-optimal transitions from multiple tasks and assigns reward labels according to the reward function. Our work is most closely related to offline RL work by Yu et al. (2022). This prior work leverages reward-free, sub-optimal data by pseudo-labeling all transitions with zero and adding them to the RL agent’s replay buffer. However, we found that directly applying this approach to the HitL RL setting was ineffective (See Appendix D.5).

3 BACKGROUND

In the RL paradigm, agents interact with an environment to maximize the total (discounted) expected reward it can achieve. This interaction process is modeled as a Markov Decision Process (MDP) which consists of $\langle \mathcal{S}, \mathcal{A}, T, r, \gamma \rangle$. At every time-step t , the agent receives a state $s_t \in \mathcal{S}$ from the environment and chooses an action $a_t \in \mathcal{A}$. The environmental transition function, T , determines the probability of transitioning to state s_{t+1} and receiving reward r_{t+1} , given the agent was in state s_t and executed action a_t . The environment then provides the agent r_{t+1} . The agent attempts to learn a policy, $\pi : \mathcal{S} \rightarrow \mathcal{A}$, that maximizes the expected return $\mathbb{E}[G] = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$, which is defined as the expected sum of discounted future rewards with discount factor $\gamma \in [0, 1)$.

3.1 REWARD LEARNING FROM HUMAN FEEDBACK

This paper assumes that we are in a reward-free paradigm, an MDP/R setting, where our goal is to (1) learn a reward function, \hat{r} , from human feedback and (2) learn a policy that maximizes the total expected \hat{r} . We follow the standard reward learning framework that uses supervised learning to learn a parameterized reward function, \hat{r}_θ , with parameters θ (Christiano et al., 2017). In both scalar- and preference-based feedback settings, we consider trajectory segments σ , where σ consists of a sequence of states and actions: $\{s_t, a_t, s_{t+1}, a_{t+1}, \dots, s_{t+k}, a_{t+k}\}$, with k as the segment size.

Preference-based Reward Learning In preference-based learning, two segments, σ^0 and σ^1 , are compared by a teacher, yielding $y \in \{0, 0.5, 1\}$. Specifically, if the teacher preferred segment σ^1 over segment σ^0 , then y is set to 1, and if the converse is true y is set to 0. If both segments are equally preferred, then y is set to 0.5. As feedback is collected, it is stored as tuples (σ^0, σ^1, y) in the reward model data set D_{RM} . In general, if $\sigma^i > \sigma^j$, then the segment σ^i is preferred by the teacher over segment σ^j . We follow the Bradley-Terry model (Bradley & Terry, 1952) to define a preference predictor using the reward function \hat{r}_θ :

$$P_\theta(\sigma^1 > \sigma^0) = \frac{\exp(\sum_t \hat{r}_\theta(s_t^1, a_t^1))}{\sum_{i \in \{0,1\}} \exp(\sum_t \hat{r}_\theta(s_t^i, a_t^i))} \quad (1)$$

Intuitively, this model assumes that the probability of the teacher preferring a segment depends exponentially on the total sum of predicted rewards along the segment. To train the reward function, we can use supervised learning where the teacher provides the labels y . More specifically, we update \hat{r}_θ by minimizing the standard binary cross-entropy objective:

$$L^{CE}(\theta, D) = - E_{(\sigma^0, \sigma^1, y) \sim D} [(1 - y) \log P_\theta(\sigma^0 > \sigma^1) + y \log P_\theta(\sigma^1 > \sigma^0)] \quad (2)$$

Scalar-based Reward Learning The primary difference between scalar and preference-based reward learning is that in scalar-based learning, the human teacher assigns numerical ratings to trajectory segments. In this setting, the comparisons between segments are implicit. More concretely, a teacher assigns a scalar value y to a segment σ^i , and as feedback is collected, it is stored as tuples (σ^i, y) in the reward model data set D_{RM} . We then apply standard regression and update \hat{r}_θ by minimizing the mean squared error:

$$L^{MSE}(\theta, D) = E_{(\sigma^i, y) \sim D} [(y - \sum_t \hat{r}_\theta(s_t^i, a_t^i))^2] \quad (3)$$

Human Studies in Human in the Loop RL To evaluate HitL RL algorithms, a common protocol is the use of simulated teachers, where feedback is provided according to a ground truth reward

function. This can be useful, as it offers an efficient and controlled evaluation setting. However, we argue that it is essential to evaluate HitL RL algorithms with *real human feedback*. This is especially important in light of recent work finding discrepancies between preference learning algorithms when using simulated teacher feedback versus human feedback (Metcalf et al., 2024).

4 SUB-OPTIMAL DATA PRE-TRAINING

In this section, we present SDP, a tool that leverages sub-optimal trajectories to improve the feedback efficiency for HitL RL. We refer to sub-optimal trajectories as sequences of (s, a) pairs such that:

$$r(s_t, a_t) - r_{\min} < \epsilon \quad \forall t \in [t, t+k] \quad \text{for some small } \epsilon > 0, \quad (4)$$

where r_{\min} is the minimum possible environment reward. Equation 4 essentially expresses that the rewards achieved along a sub-optimal trajectory should be close to r_{\min} . However, it is important to note that in practice we do not have access to the true reward. This prevents us from directly identifying sub-optimal trajectories using equation 4. Instead, we rely on selecting trajectories that we estimate will align with this criterion, such as gathering trajectories via a random policy.

Once sub-optimal trajectories are collected, we take the approach of pseudo-labeling all transitions with r_{\min} . The goal of SDP is then to use this pseudo-labeled data to create a prior for rewards models in HitL RL methods (see Figure 1). Its simplicity enables SDP to be used in conjunction with any off-the-shelf HitL RL algorithm that learns a reward function from feedback.

SDP comprises two phases: (1) the reward model pre-training phase and (2) the agent update phase. In the reward model pre-training phase, we first gather a data set, D_{sub} , of N sub-optimal state, action transitions. We then pseudo-label all transitions in D_{sub} with rewards of r_{\min} , resulting in $D_{\text{sub}} = \{s_i, a_i, r_{\min}\}_{i=1}^N$. D_{sub} is then used to optimize the reward model \hat{r}_θ with the mean squared loss in Equation 3. As a result, the reward model \hat{r}_θ learns to associate all sub-optimal transitions with a low reward. Without such a prior, the reward model would initially have random estimates for these transitions; while the only way to improve such estimates is to obtain feedback from a teacher. Therefore, the reward model pre-training phase provides a valuable reward initialization without requiring any feedback.

Next, in the agent update phase, we initialize the RL agent’s replay buffer D_{agent} with D_{sub} . The RL agent then briefly interacts with its environment and performs gradient updates according to its loss functions. The agent update process changes the RL agent’s policy and generates new transitions, which are then stored in both the agent’s replay buffer D_{agent} and the reward model’s data set D_{RM} . It is important to note that in standard scalar- and preference-based reward learning, we query the teacher for feedback on trajectory segments sampled from D_{RM} . Therefore, adding new transitions into D_{RM} during the agent update phase is necessary to ensure that the teacher does not provide redundant feedback to the original sub-optimal transitions (as D_{RM} was empty prior to the agent update phase). When it is time for the teacher to provide their first set of feedback, the feedback can cover a different region of the state and action space, relative to the original sub-optimal data. In Appendix D.1, Figure 8 we empirically show that the agent update phase changes the RL agent’s policy by performing policy rollouts and analyzing the differences in state distributions. See Algorithm 1 for the complete pseudocode.

At first glance, labeling sub-optimal transitions with an incorrect reward may seem problematic, as incorrect labels does introduce statistical bias into the reward model and the RL agent’s value network. However, as the transitions are sub-optimal, we observed that the bias for using an incorrect reward is low (see Figure 9 in Appendix D.2). Moreover, by using the sub-optimal transitions, we increase the overall amount of data used by both models, which can *decrease* the models’ variance, as shown in the offline RL setting (Yu et al., 2022).

5 EXPERIMENTS

This section considers the following four research questions (RQ’s):

RQ 1: Can SDP improve upon existing scalar- and preference-based RL methods?

RQ 2: Can SDP effectively leverage sub-optimal trajectories from different tasks to improve performance on a target task?

RQ 3: Can SDP be used with real human feedback?

RQ 4: How sensitive is SDP to various hyperparameters?

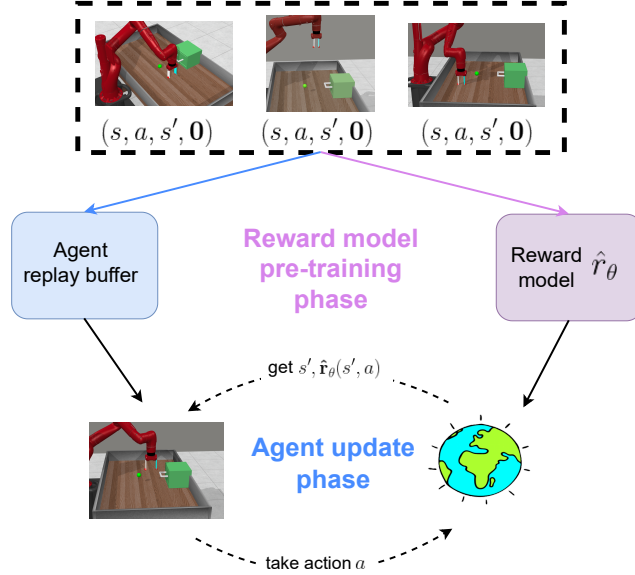


Figure 1: Overview of SDP: After obtaining sub-optimal trajectories, we pseudo-label this data with rewards of $r_{\min} = 0$. We then pre-train the reward model \hat{r}_θ using this data set. During the agent update phase, we initialize the RL agent’s replay buffer with the same pseudo-labeled data set. The agent then interacts in the environment and makes learning updates to obtain new behaviors for a teacher to give feedback.

Algorithm 1 SDP

Require: Reward model $\hat{r}_\theta \leftarrow \theta$ randomly initialized, Reward model data set $D_{\text{RM}} \leftarrow \emptyset$, RL agent with replay buffer $D_{\text{agent}} \leftarrow \emptyset$, Sub-optimal data set D_{sub} with reward labels r_{\min}

```

1: // REWARD MODEL PRE-TRAIN PHASE
2: for each gradient step do
3:     Optimize  $\hat{r}_\theta$  on  $D_{\text{sub}}$  with  $L^{\text{MSE}}$  (equ. 3)
4: end for
5: // AGENT UPDATE PHASE
6:  $D_{\text{agent}} \leftarrow D_{\text{sub}}$ 
7: for each time-step  $t$  do
8:     Collect  $s_{t+1}$  by taking action  $a_t \sim \pi(s_t)$ 
9:     Store  $(s_t, a_t, \hat{r}_\theta, s_{t+1})$  in  $D_{\text{agent}}$ 
10:    Store  $(s_t, a_t)$  in  $D_{\text{RM}}$ 
11:    Update RL agent with  $D_{\text{agent}}$ 
12: end for
13: Begin scalar- or preference-based RL using pre-trained  $\hat{r}_\theta$ , RL agent, and  $D_{\text{agent}}, D_{\text{RM}}$ 

```

5.1 EXPERIMENTAL DESIGN

To demonstrate the versatility and effectiveness of SDP, we apply SDP to both preference and scalar-based RL approaches. However, as preference feedback can be less time-consuming than scalar feedback, we primarily concentrate on preference-based RL in our experiments, exploring scalar

270 feedback in a smaller capacity. For the preference-based experiments, we combine SDP with four
 271 contemporary preference-based algorithms: PEBBLE (Lee et al., 2021a), RENE (Liang et al., 2022),
 272 SURF (Park et al., 2022), and MRN (Liu et al., 2022). We benchmark the performance of the four
 273 algorithms augmented with SDP against their original versions without SDP, as well as against
 274 SAC. We treat SAC (Haarnoja et al., 2018) as an oracle (i.e., upper bound) because it learns while
 275 accessing the true reward function, which is unavailable to the other algorithms. For the scalar-
 276 based experiments, we combine SDP with R-PEBBLE (a regression variant of PEBBLE (Lee et al.,
 277 2021a)). We compare SDP + R-PEBBLE against R-PEBBLE, Deep TAMER (Warnell et al., 2018)
 278 (a scalar feedback RL algorithm), and the (oracle) SAC. We note that SAC is the core RL algorithm
 279 used across all baselines.

280 **Implementation Details** For SDP, we collected sub-optimal trajectories via a random policy
 281 (50,000 state, action transitions for all Section 5.2 experiments). Note that we do not require ex-
 282 plicit access to a sub-optimal policy; we only require state, action transitions from said policy.
 283 Moreover, all reward model hyperparameters remained the same during both SDP phases (i.e., the
 284 reward model pre-train and agent update phases) as well as during the standard scalar- or preference-
 285 based RL that followed. As for feedback budgets, we maintained equal budgets for all algorithms
 286 within each environment. However, the budget itself was adjusted across environments to reflect
 287 their difficulty levels. Refer to Appendix A for a complete overview of the implementation process
 288 and specific hyperparameters for all algorithms.

290 **Evaluation** We show average offline performance (i.e., freeze the policy and evaluate it with no ex-
 291 ploration) over ten episodes using either the ground truth reward function (DMControl experiments)
 292 or the success rate (Meta-World experiments). It is important to note that only SAC has access to the
 293 true reward function. We perform this evaluation every 10,000 training steps. To systemically evalu-
 294 ate performance, we use a simulated teacher that provides either a scalar rating of a single trajectory
 295 segment or preferences between two trajectory segments according to the true reward function. To
 296 thoroughly test the effectiveness of SDP, we perform evaluations on four robotic locomotion tasks
 297 from the DMControl Suite: Walker-walk, Cheetah-run, Quadruped-walk, and Cartpole-swingup,
 298 and five robotic manipulation tasks from Meta-World: Hammer, Door-unlock, Door-lock, Drawer-
 299 open, and Window-open. In our experiments, the results are averaged over five seeds with shaded
 300 regions or error bars indicating 95% confidence intervals. To test for significant differences in final
 301 performance (i.e., the undiscounted return) and learning efficiency (i.e., the total area under the re-
 302 turn curve, AUC), we perform Welch t-tests (equal variances not assumed) with a p-value of 0.05.
 303 See Appendix D, Tables 9-14 for a summary of final performance and AUC across all experiments.

304 5.2 LOCOMOTION AND MANIPULATION RESULTS

306 **Preference Feedback Experiments** We first address RQ 1 by evaluating the utility of SDP in the
 307 preference-based RL setting. Considering all four preference-based algorithms in the nine environ-
 308 ments, SDP significantly ($p < 0.05$) improved learning (i.e., either final performance or AUC) in 23
 309 out of the 36 experiments (see Figure 2). In the remaining experiments, there were no significant
 310 differences in the performance between SDP and the baseline algorithms. [The addition of SDP \(i.e.,
 311 SDP + base algorithm\) never statistically hurt performance.](#)

312 **Scalar Feedback Experiments** Continuing our investigation into RQ 1, we now evaluate the
 313 performance of SDP in the scalar-based RL setting. We performed evaluations in Walker-walk,
 314 Cheetah-run, and Quadruped-walk. In Figure 3, we found that SDP (purple curve) significantly im-
 315 proves either the final performance or AUC compared to R-PEBBLE (navy curve) and Deep TAMER
 316 (yellow curve). More impressively, we found that SDP achieves comparable final performance to
 317 SAC (red curve), which uses the ground truth reward function, using as little as 60 feedback queries.
 318

319 **Leveraging Different Task Data in SDP** Our previous experiments showed that SDP can lever-
 320 age sub-optimal data from the target task to improve HitL RL methods. Now, addressing RQ 2, we
 321 investigate whether SDP can similarly use sub-optimal data from *related tasks* to improve perfor-
 322 mance on the target task. We perform three preference learning experiments, comparing PEBBLE
 323 with SDP + PEBBLE, using sub-optimal data from a different prior task that has the same virtual
 robot: (1) Walker-stand for Walker-walk, (2) Quadruped-walk for Quadruped-run, and (3) Drawer-

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

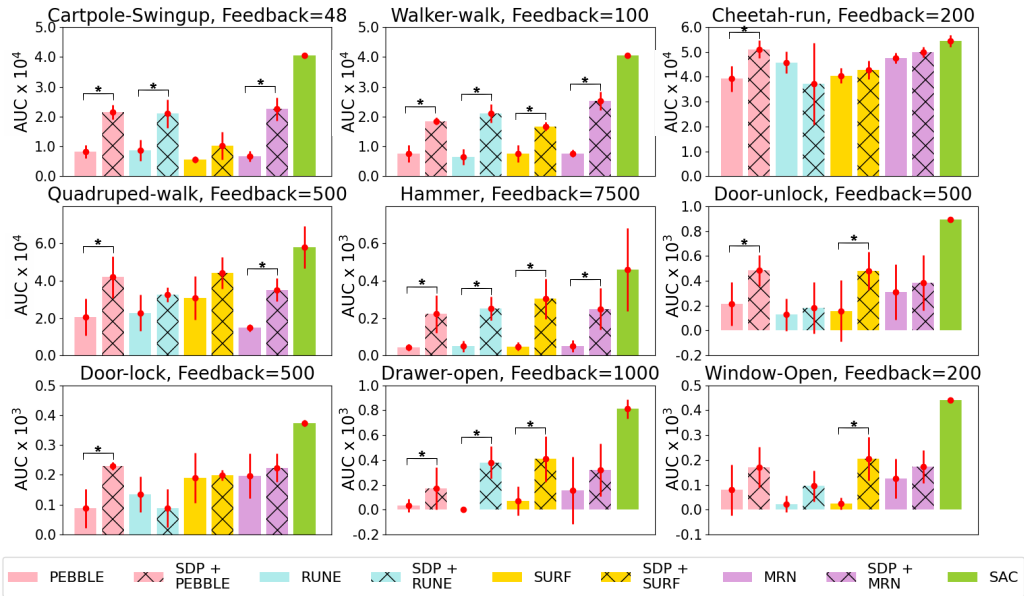


Figure 2: Preference feedback experiments in DMControl and Meta-World suites. The bar plots show AUC +/- 95% confidence intervals. * indicates that SDP + the base preference learning algorithm achieves a statistically greater score ($p < 0.05$) than the base preference learning algorithm.

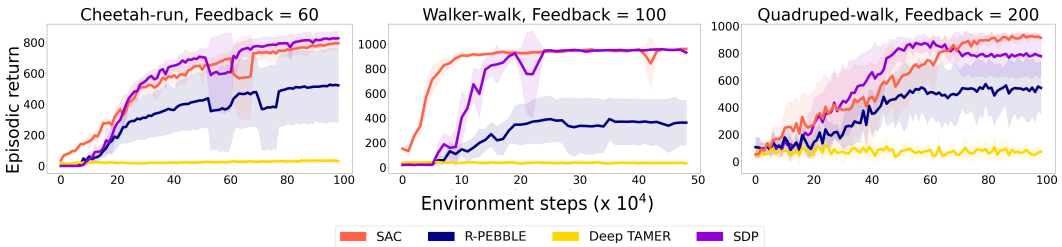


Figure 3: Scalar feedback experiments. In all DMControl experiments, SDP significantly outperforms R-PEBBLE and Deep TAMER ($p < 0.05$) and achieves comparable performance to SAC.

open for Door-open. To obtain the sub-optimal data for the prior tasks, we gathered transitions from partially trained policies as opposed to using random policies. Each partially trained policy achieved a final score of approximately 15-20% of that achieved by a fully trained policy. This ensured that the distribution of sub-optimal data differed between the prior and target tasks. See Appendix A.2 for further details on the experiment setup. Figure 4 demonstrates that in all three tested environments, SDP can successfully leverage sub-optimal data from related tasks (green curve) as it achieved similar performance to SDP when leveraging target task data (purple curve).

5.3 PREFERENCE LEARNING WITH HUMAN FEEDBACK

To evaluate HitL RL algorithms, we argue that it is critical to understand their efficacy with human teachers. However, human user studies do not appear to be widely adopted in the current literature. To empirically investigate this, we conducted a survey of 45 preference learning studies from 2012 to 2024¹ to understand the prevalence of human user studies in the preference learning literature. We found that *fewer than 50%* tested their proposed algorithms with human participants who were not also the authors. Moreover, of those studies that did involve human participants, only 41% included non-expert individuals, while none provided sufficient demographic information about their partic-

¹See Appendix B.2 for more details.

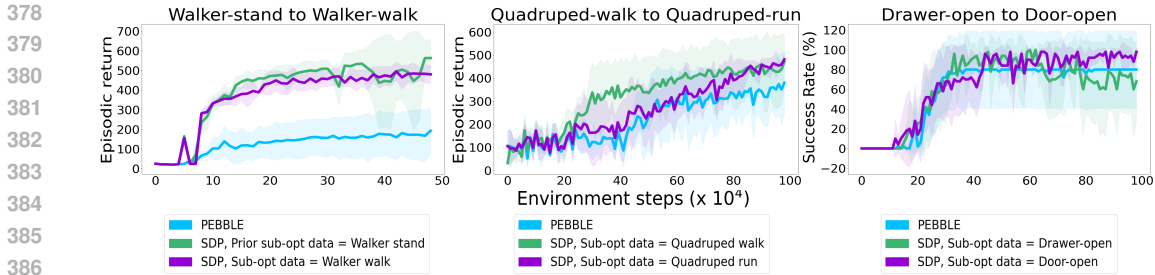


Figure 4: These figures highlight that SDP can leverage sub-optimal data from different prior tasks as it performed comparably to SDP when using target task data.

ipants (e.g., gender, race/ethnicity, level of education). These limitations raise significant concerns about whether and how current preference learning algorithms generalize to different populations.

To that end, we now address RQ 3 and evaluate the efficacy of SDP with real human feedback. We conduct an ethics-approved human-subject study of 16 participants (9 male, 7 female). Age ranges were collected: 18-24 (6), 25-30 (7), and 50-70 (3). Participants self-identified their membership in racial groups: South Asian (6), East Asian (3), White (4), Middle Eastern or North African (2), and multi-racial (1). The participants’ highest educational attainments were: high school diploma (3), Bachelor’s degrees (8), and Master’s degrees (5), and their expertise varied across AI/ML (12), other computer science topics (1), and non-computer science fields (3). This highlights the diversity of our participants in terms of demographics and expertise.

This user study compares the performance of SDP and PEBBLE in two DMControl environments: Pendulum-swingup (with 7 participants) and Cartpole-swingup (with 12 participants). We focused our comparison to PEBBLE as it performed comparably to the other preference learning baselines in Section 5.2. We use a between-subjects experimental design — each participant provides preferences for a single seed of both SDP and PEBBLE algorithms. The preference budgets for Pendulum-swingup and Cartpole-swingup were 40 and 48, respectively. We selected these environments specifically because they could be solved with fewer preferences, aiming to reduce the overall time commitment required from participants. Each trial for a single environment lasted approximately 1-1.5 hours. See Appendix B for more details including user instructions and interface.

We visualize both algorithms’ final performance and AUC in Figure 5. We found that in both environments SDP (purple plots) maintains significant ($p < 0.05$) performance gains over PEBBLE (blue plots) in terms of either final performance or AUC. In Cartpole-swingup, we also observed consistent performance from SDP regardless of whether the teacher was human or simulated. This suggests that our prior results with simulated teachers can generalize to settings where human feedback is provided. Moreover, we observed no significant differences in SDP’s effectiveness across demographic factors, including gender, age, educational, and computer science background (see Appendix C, Tables 7 and 8). This finding is particularly encouraging for the potential real-world deployment of SDP, highlighting its usability for non-expert users across diverse demographics.

5.4 ABLATION AND SENSITIVITY STUDIES

To further understand the effectiveness of SDP, we perform further analysis of SDP across three dimensions: (1) the phases of SDP, (2) the number of feedback queries, and (3) the amount of sub-optimal data. This analysis aims to address RQ 4 and provide a deeper understanding of the factors influencing SDP’s performance. For these experiments, we focus on SDP + R-PEBBLE in the Walker-walk environment. Supplementary results for Cheetah-run can be found in Appendix D.

SDP Component Analysis First, we evaluate the effect of each phase of SDP individually, the reward model pre-train phase and the agent update phase. Figure 6–leftmost demonstrates the importance of using both phases in SDP for scalar-based RL approaches. We found that the SDP variants that only use one of the phases (green and gray curves) result in worse performance than the full SDP (purple curve).

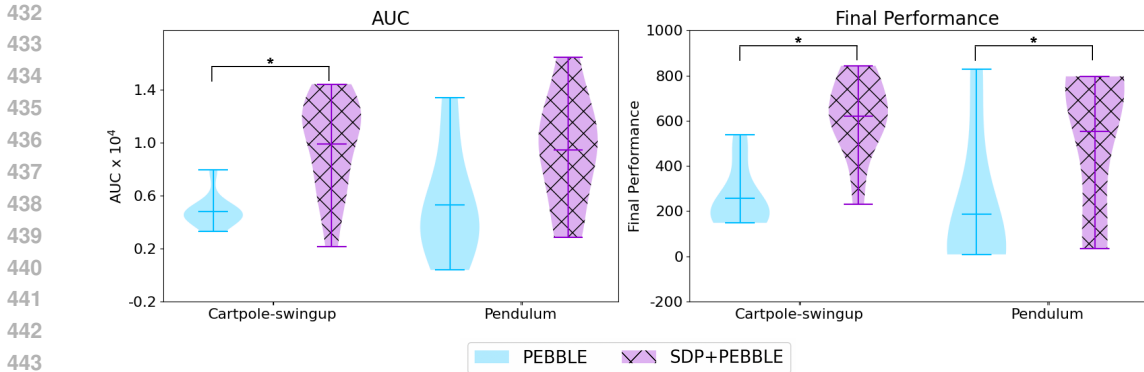


Figure 5: This figures demonstrate that SDP can significantly outperform PEBBLE in terms of both AUC (left) and final performance (right) even when human teachers are providing preferences.

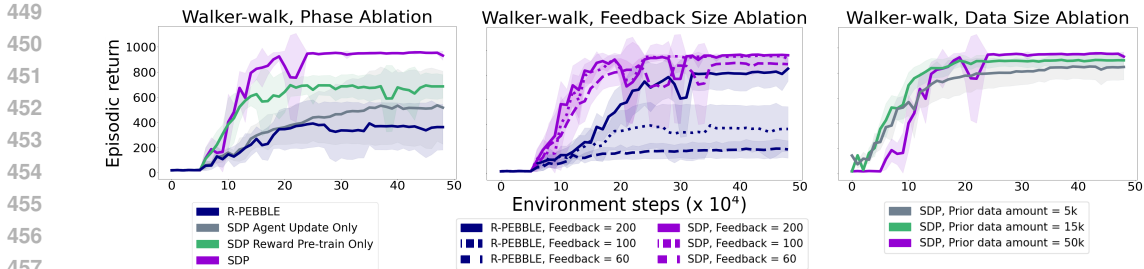


Figure 6: These figures show ablation and sensitivity studies of SDP in Walker-walk .

Effect of Feedback Amount We evaluate SDP and R-PEBBLE with feedback budgets of 60, 100, and 200 to analyze the impact of feedback quantity on performance. As shown in Figure 6 (middle), SDP (purple curves) consistently outperforms R-PEBBLE (navy curves), further demonstrating its effectiveness across varying feedback levels.

Effect of Sub-Optimal Data Amount We evaluate the performance of SDP using varying amounts of sub-optimal transitions: 5000, 15000, and 50000. Figure 6 (rightmost) indicates that while 5000 transitions (gray curve) led to the poorest performance, increasing this amount to 15000 or 50000 (green and purple curves) yielded comparable or improved results, suggesting that more sub-optimal data can benefit SDP.

6 CONCLUSION

In this work, we present SDP, an approach that improves the feedback efficiency for HitL RL algorithms. SDP is specifically designed to leverage reward-free, sub-optimal data for scalar- and preference-based HitL RL approaches. By pseudo-labeling low-quality data with the minimum environment rewards, we can pre-train the reward model without the need for human labeling. This provides the reward model with a head start in learning. This head start allows the reward model to learn to associate low-quality transitions with low reward values, even before receiving any actual human feedback. Our simulated teacher experiments in DMControl and Meta-World suites demonstrate that SDP can significantly improve a variety of preference- and scalar-based reward learning algorithms. Importantly, we further validate the real-world applicability of SDP by demonstrating its success in a 16-person user study. This work takes an important step towards considering how sub-optimal data can be leveraged for HitL RL.

ETHICS STATEMENT

The goal of preference and scalar-based RL is to learn reward functions that encode human preferences. However, this necessitates interaction with human users, raising concerns about potential biases in the collected preferences. If these preferences are primarily drawn from a non-representative group, the resulting RL system may unfairly prioritize the desires or needs of that group over others. To that end, we performed a human subject study that was reviewed and approved by an external ethics committee. This ensured the responsible and ethical collection of human preferences. In addition, we took care to gather participants from a wide range of both demographic and educational backgrounds, which we detail in Section 5.3.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we provide a link to our anonymous downloadable source code: https://anonymous.4open.science/r/SDP_ICLR-C23F/. We also provide pseudocode in Section 4. Moreover, we outline extensive training and evaluation details in Section A and Appendix A.2. We also outline all hyperparameters used for all algorithms and environments in Appendix A.2, Tables 1– 5. Lastly, for the human subject study, we provide details on the instructions provided to the participants in Appendix B as well as a description of the preference interface used. In our released codebase, we also provide the specific code files used to run the experiments in the user study. In addition, an anonymized version of the human subject data will be released, as stated in our ethics approval.

REFERENCES

- Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *Twenty-First International Conference on Machine Learning*, 2004.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *Thirty-First Conference on Neural Information Processing Systems*, 2017.
- Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 2009.
- Peter Barnett, Rachel Freedman, Justin Svegliato, and Stuart Russell. Active reward learning from multiple teachers. In *SafeAI 2023, The AAAI Workshop on Artificial Intelligence Safety*, 2023.
- Erdem Biyik and Dorsa Sadigh. Batch active preference-based learning of reward functions. In *2nd Conference on Robot Learning*, 2018. URL <http://arxiv.org/abs/1810.04303>.
- Erdem Biyik, Dylan P Losey, Malayandi Palan, Nicholas C Landolfi, Gleb Shevchuk, and Dorsa Sadigh. Learning reward functions from diverse sources of human feedback: Optimally integrating demonstrations and preferences. *The International Journal of Robotics Research*, 2022.
- Yaniv Blumenfeld, Dar Gilboa, and Daniel Soudry. Beyond signal propagation: Is feature diversity necessary in deep neural network initialization? In *Thirty-Seventh International Conference on Machine Learning*, 2020.
- Serena Booth, W. Bradley Knox, Julie Shah, Scott Niekum, Peter Stone, and Alessandro Allievi. The perils of trial-and-error reward design: Misdesign through overfitting and invalid task specifications. In *Thirty-Seventh AAAI Conference on Artificial Intelligence*, 2023.
- Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 1952.
- Daniel Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond sub-optimal demonstrations via inverse reinforcement learning from observations. In *Thirty-sixth International Conference on Machine Learning*, 2019.

- 540 Daniel S. Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning via
541 automatically-ranked demonstrations. In *Third Proceedings of the Conference on Robot Learning*,
542 2020.
- 543 Tim Brys, Anna Harutyunyan, Halit Bener Suay, Sonia Chernova, Matthew E Taylor, and Ann
544 Nowé. Reinforcement learning from demonstration through shaping. In *Twenty-Fourth Interna-*
545 *tional Joint Conference on Artificial Intelligence*, 2015.
- 546 Erdem Bıyık, Nicolas Huynh, Mykel J. Kochenderfer, and Dorsa Sadigh. Active preference-based
547 gaussian process regression for reward learning. In *Robotics: Science and Systems*, 2020. URL
548 <https://arxiv.org/abs/2005.02575>.
- 549 Serkan Cabi, Sergio Gómez Colmenarejo, Alexander Novikov, Ksenia Konyushkova, Scott Reed,
550 Rae Jeong, Konrad Zolna, Yusuf Aytar, David Budden, Mel Vecerik, et al. Scaling data-driven
551 robotics with reward sketching and batch reinforcement learning. *Robotics: Science and Systems*,
552 2020.
- 553 Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via
554 self-supervised reward regression. In *Forth Conference on Robot learning*, 2021.
- 555 Jie Cheng, Gang Xiong, Xingyuan Dai, Qinghai Miao, Yisheng Lv, and Fei-Yue Wang. RIME:
556 Robust preference-based reinforcement learning with noisy preferences. In Ruslan Salakhutdinov,
557 Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix
558 Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, vol-
559 *ume 235 of Proceedings of Machine Learning Research*, pp. 8229–8247. PMLR, 21–27 Jul 2024.
560 URL <https://proceedings.mlr.press/v235/cheng24k.html>.
- 561 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep re-
562 inforcement learning from human preferences. In *Thirty-First Conference on Neural Information*
563 *Processing Systems*, 2017.
- 564 Jack Clark and Dario Amodei. Faulty reward functions in the wild, 2016. URL <https://openai.com/research/faulty-reward-functions>.
- 565 Christian Daniel, Malte Viering, Jan Metz, Oliver Kroemer, and Jan Peters. Active reward learning.
566 In *Robotics: Science and Systems*, 2014.
- 567 Oliver Daniels-Koch and Rachel Freedman. The expertise problem: Learning from specialized
568 feedback. In *ML Safety Workshop, 36th Conference on Neural Information Processing Systems*
569 *(NeurIPS 2022)*, 2022. URL <https://arxiv.org/abs/2211.06519>.
- 570 Anca Dragan and Siddhartha Srinivasa. Formalizing assistive teleoperation. *Robotics: Science and*
571 *Systems*, 2012.
- 572 Yang Gao, Huazhe Xu, Ji Lin, Fisher Yu, Sergey Levine, and Trevor Darrell. Reinforcement learning
573 from imperfect demonstrations. In *Thirty-Fifth International Conference on Machine Learning*,
574 2019.
- 575 Joseph Giovanelli, Alexander Tornede, Tanja Tornede, and Marius Lindauer. Interactive hyperpa-
576 rameter optimization in multi-objective problems via preference learning. In *Proceedings of the*
577 *AAAI Conference on Artificial Intelligence*, 2024.
- 578 Shane Griffith, Kaushik Subramanian, Jonathan Scholz, Charles L Isbell, and Andrea L Thomaz.
579 Policy shaping: Integrating human feedback with reinforcement learning. In *Twenty-Sixth Con-*
580 *ference on Neural Information Processing Systems*, 2013.
- 581 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
582 maximum entropy deep reinforcement learning with a stochastic actor. In *Thirty-Fifth Interna-*
583 *tional Conference on Machine Learning*, 2018.
- 584 Donald Joseph Hejna III and Dorsa Sadigh. Few-shot preference learning for human-in-the-loop rl.
585 In *Sixth Conference on Robot Learning*, 2023.

- 594 Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
595 John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations. In
596 *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
597
- 598 Simon Holk, Daniel Marta, and Iolanda Leite. Predilect: Preferences delineated with zero-shot
599 language-based reasoning in reinforcement learning. In *Proceedings of the 2024 ACM/IEEE In-*
600 *ternational Conference on Human-Robot Interaction*, volume 32 of *HRI '24*, pp. 259–268. ACM,
601 March 2024. doi: 10.1145/3610977.3634970. URL [http://dx.doi.org/10.1145/](http://dx.doi.org/10.1145/3610977.3634970)
602 [3610977.3634970](http://dx.doi.org/10.1145/3610977.3634970).
- 603 Xiao Hu, Jianxiong Li, Xianyuan Zhan, Qing-Shan Jia, and Ya-Qin Zhang. Query-policy mis-
604 alignment in preference-based reinforcement learning. In *The Twelfth International Confer-*
605 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=UoBymIwPJr)
606 [UoBymIwPJr](https://openreview.net/forum?id=UoBymIwPJr).
- 607
- 608 Minyoung Hwang, Gunmin Lee, Hogun Kee, Chan Woo Kim, Kyungjae Lee, and Songhwai Oh.
609 Sequential preference ranking for efficient reinforcement learning from human feedback. In
610 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL [https://](https://openreview.net/forum?id=MIYBTjCVjR)
611 openreview.net/forum?id=MIYBTjCVjR.
- 612 Minyoung Hwang, Luca Weihs, Chanwoo Park, Kimin Lee, Aniruddha Kembhavi, and Kiana
613 Ehsani. Promptable behaviors: Personalizing multi-objective rewards from human pref-
614 erences. In *Conference on Computer Vision and Pattern Recognition*, 2023b. URL
615 [https://openaccess.thecvf.com/content/CVPR2024/papers/Hwang_](https://openaccess.thecvf.com/content/CVPR2024/papers/Hwang_Promptable_Behaviors_Personalizing_Multi-Objective_Rewards_from_Human_Preferences_CVPR_2024_paper.pdf)
616 [Promptable_Behaviors_Personalizing_Multi-Objective_Rewards_from_](https://openaccess.thecvf.com/content/CVPR2024/papers/Hwang_Promptable_Behaviors_Personalizing_Multi-Objective_Rewards_from_Human_Preferences_CVPR_2024_paper.pdf)
617 [Human_Preferences_CVPR_2024_paper.pdf](https://openaccess.thecvf.com/content/CVPR2024/papers/Hwang_Promptable_Behaviors_Personalizing_Multi-Objective_Rewards_from_Human_Preferences_CVPR_2024_paper.pdf).
- 618
- 619 Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward
620 learning from human preferences and demonstrations in atari. In *Thirty-Second Conference on*
621 *Neural Information Processing Systems*, 2018.
- 622 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Third Interna-*
623 *tional Conference on Learning Representations*, 2015.
624
- 625 W Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement. In *Fifth*
626 *International Conference on Knowledge Capture*, 2009.
- 627
- 628 W Bradley Knox and Peter Stone. Learning non-myopically from human-generated reward. In
629 *International Conference on Intelligent User Interfaces*, 2013.
- 630 W. Bradley Knox, Stephane Hatgis-Kessell, Serena Booth, Scott Niekum, Peter Stone, and Alessan-
631 dro Allievi. Models of human preference for learning reward functions. *Transactions on Machine*
632 *Learning Research (TMLR)*, 2023.
- 633
- 634 Kimin Lee, Laura Smith, and Pieter Abbeel. Pebble: Feedback-efficient interactive reinforcement
635 learning via relabeling experience and unsupervised pre-training. In *Thiry-Eighth International*
636 *Conference on Machine Learning*, 2021a.
- 637
- 638 Kimin Lee, Laura Smith, Anca Dragan, and Pieter Abbeel. B-pref: Benchmarking preference-
639 based reinforcement learning. In *Thirty-fifth Conference on Neural Information Processing Sys-*
640 *tems Datasets and Benchmarks Track (Round 1)*, 2021b. URL [https://openreview.net/](https://openreview.net/forum?id=ps95-mkHF_)
641 [forum?id=ps95-mkHF_](https://openreview.net/forum?id=ps95-mkHF_).
- 642
- 643 Xinran Liang, Katherine Shu, Kimin Lee, and Pieter Abbeel. Reward uncertainty for exploration in
644 preference-based reinforcement learning. In *Tenth International Conference on Learning Repre-*
645 *sentations*, 2022.
- 646
- 647 Runze Liu, Fengshuo Bai, Yali Du, and Yaodong Yang. Meta-reward-net: Implicitly differentiable
648 reward learning for preference-based reinforcement learning. In Alice H. Oh, Alekh Agarwal,
649 Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Sys-*
650 *tems*, 2022. URL <https://openreview.net/forum?id=OZKBRUF-wX>.

- 648 Runze Liu, Yali Du, Fengshuo Bai, Jiafei Lyu, and Xiu Li. PEARL: Zero-shot cross-task preference
649 alignment and robust reward learning for robotic manipulation. In Ruslan Salakhutdinov, Zico
650 Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp
651 (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of
652 *Proceedings of Machine Learning Research*, pp. 30946–30964. PMLR, 21–27 Jul 2024. URL
653 <https://proceedings.mlr.press/v235/liu24o.html>.
- 654 Yi Liu, Gaurav Datta, Ellen Novoseller, and Daniel S. Brown. Efficient preference-based rein-
655 forcement learning using learned dynamics models. In *2023 IEEE International Conference on*
656 *Robotics and Automation (ICRA)*, 2023. URL <https://arxiv.org/abs/2301.04741>.
- 657 Robert Loftin, Bei Peng, James MacGlashan, Michael L Littman, Matthew E Taylor, Jeff Huang, and
658 David L Roberts. Learning behaviors via human-delivered discrete feedback: modeling implicit
659 feedback strategies to speed up learning. *Autonomous Agents and Multi-Agent Systems*, 2016.
- 660 James MacGlashan, Mark K Ho, Robert Loftin, Bei Peng, David Roberts, Matthew E Taylor, and
661 Michael L Littman. Interactive learning from policy-dependent human feedback. In *Thirty-Forth*
662 *International Conference on Machine Learning*, 2017.
- 663 Daniel Marta, Simon Holk, Christian Pek, Jana Tumova, and Iolanda Leite. Variquery: Vae segment-
664 based active learning for query selection in preference-based reinforcement learning. In *2023*
665 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 7878–7885,
666 2023a. doi: 10.1109/IROS55552.2023.10341795.
- 667 Daniel Marta, Simon Holk, Christian Pek, Jana Tumova, and Iolanda Leite. Aligning human prefer-
668 ences with baseline objectives in reinforcement learning. In *2023 IEEE International Conference*
669 *on Robotics and Automation (ICRA)*, pp. 7562–7568. IEEE, 2023b.
- 670 Daniel Marta, Simon Holk, Christian Pek, and Iolanda Leite. Sequel: Semi-supervised preference-
671 based rl with query synthesis via latent interpolation. In *2024 IEEE International Conference on*
672 *Robotics and Automation (ICRA)*, pp. 9585–9592. IEEE, 2024.
- 673 Thommen George Karimpanal Maxence Hussonnois and Santu Rana. Controlled diversity with pref-
674 erence: Towards learning a diverse set of desired skills. In *Proceedings of the 2023 International*
675 *Conference on Autonomous Agents and Multiagent Systems*, 2023.
- 676 Shaunak A Mehta and Dylan P Losey. Unified learning from demonstrations, corrections, and prefer-
677 ences during physical human–robot interaction. *ACM Transactions on Human-Robot Interaction*,
678 13(3):1–25, 2024.
- 679 Katherine Metcalf, Miguel Sarabia, Natalie Mackraz, and Barry-John Theobald. Sample-efficient
680 preference-based reinforcement learning with dynamics aware rewards. In *7th Annual Conference*
681 *on Robot Learning*, 2023. URL <https://openreview.net/forum?id=i84V7i6KEMd>.
- 682 Katherine Metcalf, Miguel Sarabia, Masha Fedzechkina, and Barry-John Theobald. Can you rely on
683 synthetic labellers in preference-based reinforcement learning? it’s complicated. In *Proceedings*
684 *of the AAAI Conference on Artificial Intelligence*, 2024. doi: 10.1609/aaai.v38i9.28877. URL
685 <https://ojs.aaai.org/index.php/AAAI/article/view/28877>.
- 686 Vivek Myers, Erdem Biyik, Nima Anari, and Dorsa Sadigh. Learning multimodal rewards from
687 rankings. In *Fifth Conference on Robot learning*, 2022.
- 688 Vivek Myers, Erdem Biyik, and Dorsa Sadigh. Active reward learning from online preferences. In
689 *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 7511–7518, 2023.
690 URL <https://api.semanticscholar.org/CorpusID:257219724>.
- 691 Jongjin Park, Younggyo Seo, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Surf:
692 Semi-supervised reward learning with data augmentation for feedback-efficient preference-based
693 reinforcement learning. In *Tenth International Conference on Learning Representations*, 2022.
- 694 Zhizhou Ren, Anji Liu, Yitao Liang, Jian Peng, and Jianzhu Ma. Efficient meta reinforcement
695 learning for preference-based fast adaptation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
696 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL
697 <https://openreview.net/forum?id=6lUwgeIotn>.

- 702 Kyriacos Shiarlis, Joao Messias, and Shimon Whiteson. Inverse reinforcement learning from failure.
703 In *Fifteenth International Conference on Autonomous Agents and Multiagent Systems*, 2016.
704
- 705 Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. Cog: Con-
706 necting new skills to past experience with offline reinforcement learning. In *Conference on Robot*
707 *Learning*, 2020.
- 708 Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and charac-
709 terizing reward gaming. In *Thirty-Sixth Conference on Neural Information Processing Systems*,
710 2022.
- 711 Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
712
- 713 Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A mini-
714 maximalist approach to reinforcement learning from human feedback. In *Proceedings of the 41st*
715 *International Conference on Machine Learning*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2401.04056)
716 [2401.04056](https://arxiv.org/abs/2401.04056).
- 717 Voot Tangkaratt, Nontawat Charoenphakdee, and Masashi Sugiyama. Robust imitation learning
718 from noisy demonstrations. In *Twenty-Forth International Conference on Artificial Intelligence*
719 *and Statistics*, 2021.
720
- 721 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud-
722 den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv*
723 *preprint arXiv:1801.00690*, 2018.
- 724 Matthew E Taylor, Halit Bener Suay, and Sonia Chernova. Integrating reinforcement learning with
725 human demonstrations of varying ability. In *Tenth International Conference on Autonomous*
726 *Agents and Multiagent Systems*, 2011.
727
- 728 Mudit Verma and Katherine Metcalf. Hindsight PRIORs for reward learning from human pref-
729 erences. In *The Twelfth International Conference on Learning Representations*, 2024. URL
730 <https://openreview.net/forum?id=NLevOah0CJ>.
- 731 Mudit Verma, Siddhant Bhambri, and Subbarao Kambhampati. Exploiting unlabeled data for
732 feedback efficient human preference based reinforcement learning. In *The AAAI 2023 Work-*
733 *shop on Representation Learning for Responsible Human-Centric AI*, 2023. URL [https:](https://arxiv.org/abs/2302.08738)
734 [//arxiv.org/abs/2302.08738](https://arxiv.org/abs/2302.08738).
- 735 Ngo Anh Vien, Wolfgang Ertel, and Tae Choong Chung. Learning via human feedback in continuous
736 state and action spaces. *Applied intelligence*, 2013.
737
- 738 Ren-Jian Wang, Ke Xue, Yutong Wang, Peng Yang, Haobo Fu, Qiang Fu, and Chao Qian. Diver-
739 sity from human feedback. In *2nd Workshop on Agent Learning in Open-Endedness (ALOE) at*
740 *NeurIPS 2023*, 2023. URL <https://arxiv.org/abs/2310.06648>.
- 741 Ruiqi Wang, Weizheng Wang, and Byung-Cheol Min. Feedback-efficient active preference learning
742 for socially aware robot navigation. In *2022 IEEE/RSJ International Conference on Intelligent*
743 *Robots and Systems (IROS)*, pp. 11336–11343. IEEE, 2022.
744
- 745 Xiaofei Wang, Kimin Lee, Kourosh Hakhmaneshi, Pieter Abbeel, and Michael Laskin. Skill pref-
746 erences: Learning to extract and execute robotic skills from human feedback. In *Fifth Conference*
747 *on Robot Learning*, 2021.
- 748 Garrett Warnell, Nicholas Waytowich, Vernon Lawhern, and Peter Stone. Deep tamer: Interactive
749 agent shaping in high-dimensional state spaces. In *Thirty-Second AAAI Conference on Artificial*
750 *Intelligence*, 2018.
- 751 Devin White, Mingkan Wu, Ellen Novoseller, Vernon J Lawhern, Nicholas Waytowich, and Yong-
752 can Cao. Rating-based reinforcement learning. In *Proceedings of the AAAI Conference on Artifi-*
753 *cial Intelligence*, 2024.
754
- 755 Nils Wilde, Erdem Biyik, Dorsa Sadigh, and Stephen L Smith. Learning reward functions from
scale feedback. In *Conference on Robot Learning*, 2021.

- 756 Aaron Wilson, Alan Fern, and Prasad Tadepalli. A bayesian approach for policy learning from
757 trajectory preference queries. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger
758 (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.,
759 2012. URL [https://proceedings.neurips.cc/paper_files/paper/2012/
760 file/16c222aa19898e5058938167c8ab6c57-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/16c222aa19898e5058938167c8ab6c57-Paper.pdf).
- 761 Wanqi Xue, Bo An, Shuicheng Yan, and Zhongwen Xu. Reinforcement learning from diverse hu-
762 man preferences. In *Proceedings of the Thirty-Third International Joint Conference on Artificial
763 Intelligence (IJCAI-24)*, 2024.
- 764 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
765 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
766 In *Third Conference on Robot Learning*, 2020.
- 767 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn.
768 Conservative data sharing for multi-task offline reinforcement learning. *Advances in Neural In-
769 formation Processing Systems*, 2021.
- 770 Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Chelsea Finn, and Sergey Levine.
771 How to leverage unlabeled data in offline reinforcement learning. In *Thirty-Ninth International
772 Conference on Machine Learning*, 2022.
- 773 Guoxi Zhang and Hisashi Kashima. Learning state importance for preference-based reinforcement
774 learning. *Machine Learning*, 113(4):1885–1901, 2024.
- 775 Jiawei Zhao, Florian Tobias Schaefer, and Anima Anandkumar. Zero initialization: Initializing
776 neural networks with only zeros and ones. *Transactions on Machine Learning Research*, 2022.
- 777 Tianchen Zhu, Yue Qiu, Haoyi Zhou, and Jianxin Li. Decoding global preferences: Temporal and
778 cooperative dependency modeling in multi-agent preference-based reinforcement learning. In
779 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. doi: 10.1609/aaai.v38i15.
780 29666. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29666>.
- 781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 APPENDIX

811
812
813 A SIMULATED EXPERIMENT DETAILS

814
815 A.1 BENCHMARKS

816
817 **PEBBLE** PEBBLE has two primary components: unsupervised exploration and off-policy learn-
818 ing with relabeling. The purpose of the unsupervised exploration phase is to collect diverse expe-
819 riences for a human teacher to provide feedback on. More specifically, PEBBLE optimizes state
820 entropy to explore the environment. Furthermore, PEBBLE uses off-policy reinforcement learning
821 to learn a policy. PEBBLE specifically uses an off-policy RL algorithm as they are more sample
822 efficient compared to their on-policy counterparts. Then as the reward model changes, PEBBLE
823 relabels all transitions in the RL agent’s replay buffer with the latest reward model. This is integral
824 as the reward model is non-stationary, and relabeling the transitions stabilizes the learning process.

825 To adapt PEBBLE to the scalar feedback setting, we make one minor change to the reward model.
826 In the scalar feedback setting, we use a scripted teacher that provides a scalar rating of a single
827 trajectory segment. Therefore, the only update to PEBBLE is with respect to the loss function.
828 Instead of using the cross-entropy loss in Equation 2, we use the mean-squared error loss in Equation
829 3.

830
831 **RUNE** RUNE is a preference learning algorithm (built on top of PEBBLE) that uses an
832 uncertainty-based exploration strategy to improve feedback efficiency. To encourage exploration
833 for the SAC agent, RUNE adds an intrinsic reward component based on the standard deviation in
834 the reward model.

835
836
837 **SURF** SURF is another preference learning algorithm (built on top of PEBBLE) that improves
838 feedback efficiency by using semi-supervised and data augmentation approaches. To incorporate
839 semi-supervised learning, SURF generates pseudo-labels for unlabeled trajectories by querying the
840 learned reward model. If the reward model confidently (e.g., low output standard deviation) pre-
841 dicted the pseudo-label, then the trajectory, label pair is added to the reward model training data set.
842 Further, SURF proposes a new data augmentation technique that crops sub-sequences of trajectories.

843
844 **MRN** MRN is a preference learning algorithm also integrated with PEBBLE. Unlike other PbRL
845 algorithms, MRN is a bi-level optimization algorithm in which the actor and critic are updated in the
846 inner loop, and the reward model is updated in the outer loop. Importantly, the reward model takes
847 into account the performance of the critic on the preference data.

848
849 **Deep TAMER** In our scalar feedback experiments, we also consider the Deep TAMER bench-
850 mark. In Deep TAMER, scalar feedback is used to learn a human reward function via regression.
851 Then the agent acts greedily with respect to this reward function. Furthermore, the original im-
852 plementation of Deep TAMER was built on top of DQN. Therefore, there was no separate actor-
853 network. In addition, Deep TAMER only used discrete feedback $\in [-1, 0, 1]$.

854 To make Deep TAMER a fair benchmark, we made a few adjustments. To start, we allow Deep
855 TAMER to learn from real-valued feedback as done in the other scalar-based experiments. However,
856 instead of using the ground truth reward function as feedback, we use the state-action values from a
857 fully-trained SAC agent. We do this because, in TAMER (and Deep TAMER), the teacher is intended
858 to provide feedback representative of the return. Secondly, in Deep TAMER, feedback is provided
859 per (state, action) pair. Therefore, to make sure Deep TAMER received the same amount of feedback
860 as the other baselines, we used *trajectory segment size* \times *feedback amount* for Deep TAMER only.
861 The other benchmarks receive a scalar feedback value that is the sum of rewards along a trajectory
862 segment. Third, to learn the reward model we use standard regression as described in Section 3.1.
863 Lastly, as our testing environments are continuous state and action, we learn a separate actor policy,
similarly done in (Vien et al., 2013).

A.2 TRAINING DETAILS

In all of our experiments, we use the hyperparameters in Table 1 for the reward models used in all benchmarks. For the agent update phase of SDP, an additional hyperparameter is associated with the number of environment interactions made before the standard preference/scalar feedback learning loop begins. However, for simplicity, we kept the same value as the existing feedback frequency hyperparameter, as feedback frequency also dictates the number of environment interactions made between feedback sessions

Furthermore, we use most of the existing reward model hyperparameters used in PEBBLE, however, we adjusted the following four hyperparameters: feedback frequency, amount of feedback per session, trajectory segment size (only for Meta-world), and activation function for the final NN layer. We adjusted the first two hyperparameters because PEBBLE originally used a significantly larger feedback budget, therefore we wanted the feedback schedule to better reflect a smaller feedback budget. We used a different trajectory segment size for Meta-world because we wanted to keep the segment sizes the same across both the DMControl and Meta-world environments. Moreover, we found that the output activation function could significantly affect learning, therefore we tested all benchmarks using both Tanh (original activation used) and Leaky-ReLU and chose the reward model that achieved better final performance. For the RUNE and SURF baselines, we use any hyperparameters associated with their specific algorithm according to the original paper (see Table 2). For a fair comparison with SDP, we provide all HiTL baselines (e.g., PEBBLE, R-PEBBLE, Deep TAMER, RUNE, and SURF) with the sub-optimal data set to be used in both the reward model and by the RL agent.

Furthermore, to select trajectory segments for the teacher to provide feedback on, we use uniform sampling in the DMControl tasks and disagreement sampling in the Meta-world tasks. Disagreement sampling is a popular active learning approach in which trajectories with higher uncertainty (based on an ensemble of neural networks) are more likely to be sampled (Christiano et al., 2017). As for the SAC hyperparameters, we use the values found in Tables 3-4.

For the experiments in which we leveraged sub-optimal data from a different task (i.e., Walker-stand, Quadruped-walk, Drawer-open), we gathered 50,000 transitions from partially trained policies. We note that for these experiments, we purposely did not use transitions gathered from a random policy. In these experiments, the prior and target tasks were environments in which the simulated robot was identical. The only difference is the environmental reward. Therefore, the random policy for both environments would be the same. To truly demonstrate transfer, we wanted to ensure we obtained low-quality transitions of the prior task that were different from the target task.

Each partially trained policy achieved a final score of approximately 15-20% of that achieved by a fully trained policy. More specifically, we used the following procedure to train the SAC policies. First, for Walker-stand, we trained a SAC policy for 5,000 time steps, and the average final performance was approximately 194. Second, for Quadruped-walk, we trained a SAC policy for 100,000 timesteps, and the average final performance was approximately 184. Lastly, for Drawer-open, we trained a SAC policy for 50,000 time steps, and the average final success rate was approximately 14%.

Table 1: Hyperparameters for the reward model used in all experiments (both preference and scalar feedback variants).

Hyperparameter	Value
segment size	50
number of random steps (i.e., sub-optimal data transitions)	50000
number of unsupervised exploration steps	9000
consider unsupervised explorations as sub-optimal	False (Walker-walk, Cartpole-swingup, Quadruped-walk – PbRL) True (others)
frequency of feedback	20000 (DMControl) 10000 (Window-open, Door-unlock, Door-lock) 5000 (Hammer, Drawer-open)
feedback budgets in ablations	2000 (Door-open) 1000 (Cheetah-run, Walker-walk)
number of feedback queries per session	50 (Hammer, Drawer-open) 8 (Cartpole-swingup) 20 (others)
sampling scheme	disagreement sampling (Metaworld) uniform sampling (DMControl)
number of training epochs	200 (Window-open, SURF and SDP + SURF) 50 (others)
learning rate	3×10^{-4}
intermediate neural network activation	Leaky ReLU
batch size	128
number of hidden layers	4
number of neurons per hidden layer	128
loss function	Mean Squared error (Scalar feedback) Cross Entropy loss (Preference feedback)
optimizer	Adam (Kingma & Ba, 2015)
<i>For Reward Model Pre-training Phase in SDP</i>	
last layer neural network activation	Tanh
<i>For Human-in-the-Loop Alg</i>	
last layer neural network activation	Leaky ReLU (PEBBLE, RUNE, SURF) Tanh (SDP – Hammer, Drawer-open, Walker-walk) Leaky ReLU (SDP – others)

Table 2: Baseline Specific Hyperparameters

Hyperparameter	Value
<i>Specific RUNE Hyperparameters</i>	
beta schedule	linear decay
beta init	0.05
beta decay	10^{-5}
<i>Specific SURF Hyperparameters</i>	
threshold λ	0.99
unlabeled batch ratio	4
loss weight	1
min/max length of cropped segment	[45, 55]
segment length before cropping	60
<i>Specific MRN Hyperparameters</i>	
num meta steps	1000 (Walker-walk) 3000(Quadruped-walk) 10000 (Door-open) 5000 (others)

972 Table 3: Hyperparameters for SAC that were shared by all algorithms. SAC is the standard RL
 973 algorithm that learns from the environment’s true reward signal, not via preference learning.
 974

Hyperparameter	Value
optimizer	Adam (Kingma & Ba, 2015)
discount	0.99
alpha learning rate	10^{-4}
actor betas	0.9, 0.999
critic betas	0.9, 0.999
alpha betas	0.9, 0.999
target smoothing coefficient	0.005
actor update frequency	1
critic target update frequency	2
init temperature	0.1
network type	MLP
nonlinearity	ReLU
number of gradient updates per step	1

985 Table 4: Specific SAC hyperparameters that were tuned for the DMControl and Metaworld experi-
 986 ments. The majority of these hyperparameters were selected from PEBBLE repo (Lee et al., 2021a).
 987

Hyperparameter	Value
<i>DMControl</i>	
batch size	512 (Cartpole-swingup) 1024 (others)
number of hidden layers	2
number of neurons per hidden layer	256 (Cartpole-swingup) 1024 (others)
actor/critic learning rate	5×10^{-5} (Cheetah-run preference feedback) 10^{-4} (Cheetah-run scalar feedback) 5×10^{-4} (Walker-walk) 10^{-4} (others)
number of training steps	0.5×10^6 (Walker-walk, Cartpole-swingup) 10^6 (others)
<i>Metaworld</i>	
batch size	512
number of hidden layers	3
number of neurons per hidden layer	256
actor/critic learning rate	3×10^{-4}
number of training steps	0.5×10^6 (Door-lock and Window-open) 10^6 (others)

1009 Table 5: Tuned hyperparameters for human feedback experiments.
 1010

Hyperparameter	Value
<i>SAC</i>	
number of training steps	0.3×10^6
batch size	1024
number of hidden layers	2
number of neurons per hidden layer	1024
actor/critic learning rate	10^{-4}
number of training steps	2.5×10^5
<i>Reward Model</i>	
consider unsupervised explorations as sub-optimal	False
frequency of feedback	20000
number of feedback queries per session	10 (Pendulum-swingup) 8 (Cartpole-swingup)
sampling scheme	uniform sampling
batch size	128
last layer neural network activation	Leaky ReLU

1026 B HUMAN SUBJECT STUDY DETAILS

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

B.1 PROTOCOL STEPS

1038

1039

1040

1041

1042

1043

The study took place in person, where the participant was first provided the consent from to review and sign. Then, participants and the researcher, together, reviewed the instructions outlining the objective of the study. The instructions had two components. First, they outlined participant’s goal in the study. Second, as done in Christiano et al. (2017), we provided a guide on what constitutes good and bad behavior in each domain.

1044

1045

Pendulum-swingup The first set of instructions for the study are as follows:

1046

1047

1048

1049

1050

1051

1052

1. Objective: for the pole to swing up and balance.
2. Agent controls how much torque (or force or twist) to apply
3. It does not matter which way the pole swings, left or right, as long as the pole balances.
4. Your task: You will see two clips of behaviors, and your job is to select the clip you think is better given the above objective. If you think both clips are identical in behavior, you can select equally preferable.

1053

1054

The second set of instructions (i.e., advice) for the study are as follows:

1055

1056

1057

1058

1059

1060

1061

1. The first priority is for the pole to begin swinging back and forth (i.e., picking up momentum). Therefore, the video clip where the pole has swung higher is better. Even if the agent is not behaving well in either clips, if you can tell that the pole is higher in one clip than the other, it is better to prefer that clip.
2. In general, if the pole is swinging rapidly in a circle (e.g., complete 360), then this is usually worse behavior than if the pole is barely moving.

1062

1063

1064

1065

1066

1067

1068

1069

1070

Cartpole-swingup The first set of instructions for the study are as follows:

1071

1072

1073

1074

1075

1076

1077

1078

1079

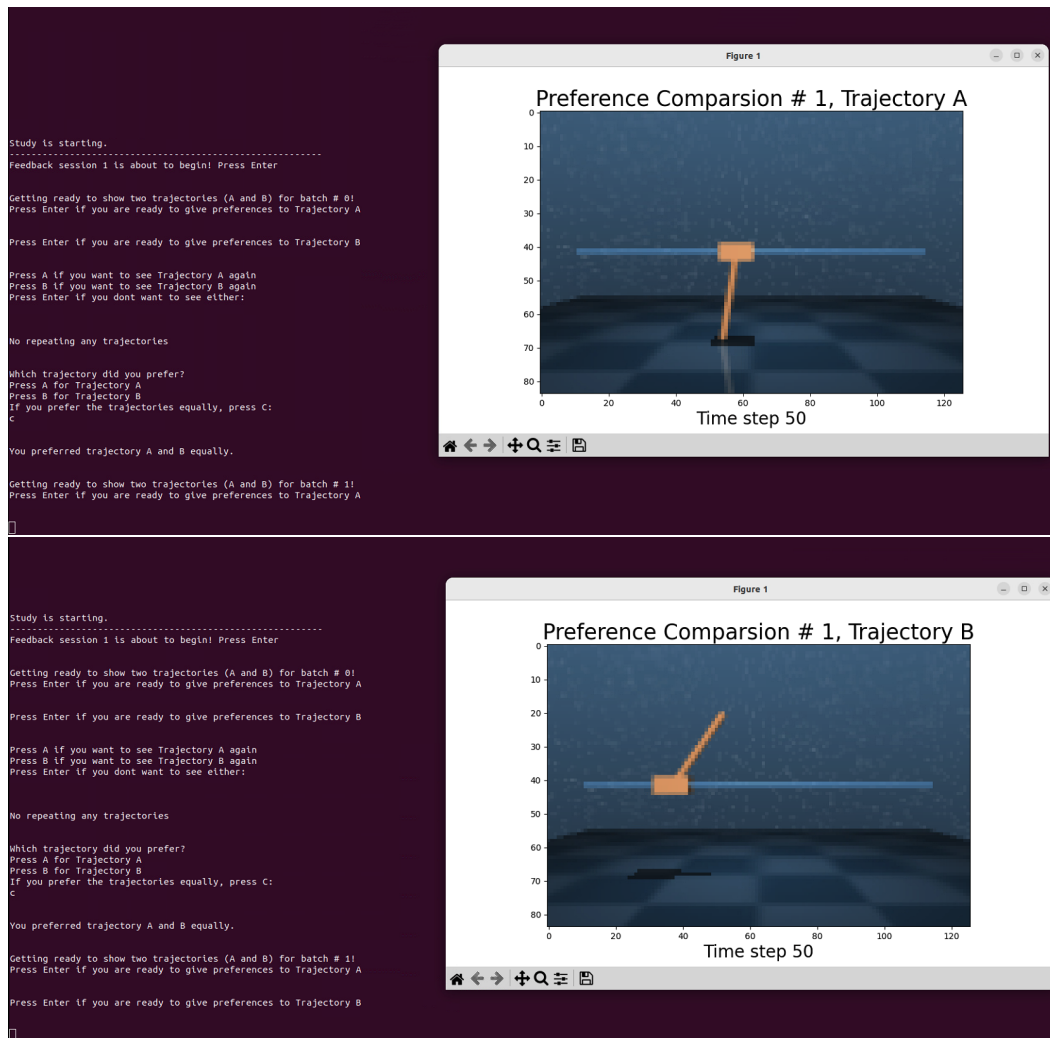
1. Pole is attached to a cart.
2. Agent can move the cart left and right.
3. Goal is for the agent to move the cart such that the pole swings up and balances
4. Your task: You will see two clips of behaviors, and your job is to select the clip you think is better given the above objective. If you think both clips are identical in behavior, you can select equally preferable.

The second set of instructions (i.e., advice) for the study are as follows:

1. The first priority is for the pole to begin swinging back and forth (i.e., picking up momentum). Therefore, the video clip where the pole has swung higher is better. Even if the agent is not behaving well in either clips, if you can tell that the pole is higher in one clip than the other, it is better to prefer that clip.
2. In general, if the pole is swinging rapidly in a circle (e.g., complete 360), then this is usually worse behavior than if the pole is barely moving.
3. If the pole is balancing then falls over, then is is better behavior than if the pole is barely moving.

1080 During the instruction period, participants also watched video clips of what successful and unsuccessful behavior looks like. In addition, participants were able to ask questions about the environment’s objective.

1083 Once the instruction period was complete, participants completed a practice round in providing preferences. This was included so users could gain familiarity with the interface. Figure 7 shows a screenshot of the interface used. Participants would use keyboard input to move from one video clip to the next. Then, participants had the opportunity to rewatch either clips as many times as they liked. Afterwards, participants chose which clip they preferred (or equally preferred) via keyboard input. Each video clip consisted of a 50 step segment, which was approximately 2 seconds long. The total study time was $\sim (1 - 1.5)$ hours. After the study was complete, participants filled out a demographic survey which included questions pertaining to their age, race/ethnicity, education level completed, and area of expertise.



1126 Figure 7: This shows the user interface used for the human subject study. Participants would view
1127 each video clip, sequentially, then decide which clip they preferred.

1128
1129
1130
1131
1132
1133

1134 B.2 PREFERENCE LEARNING SURVEY DETAILS
1135

1136 Table 6 provides a list of all papers included in our preference learning survey. We now describe
1137 the criteria for inclusion and exclusion in our survey. To find papers, we used Google Search with
1138 key words: preference learning, reinforcement learning from human feedback, and preference based
1139 reinforcement learning. We also found papers by reviewing those that cited popular preference
1140 learning works. This included Christiano et al. (2017) and Lee et al. (2021a), as these were two of
1141 the first papers that popularized preference learning. We included conference, journal and workshop
1142 papers in our survey. We did not include papers that were strictly on Arxiv. To keep the survey within
1143 scope, we did not include any preference learning works related to large language or foundation
1144 models. We also did not include any works that were entirely theoretical (e.g., no empirical results
1145 were provided).

1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

Table 6: Articles used in Preference Learning Survey from Section 5.3. NA indicates no information was provided.

Paper Reference	Human Users(Non-Author)	Number of Users
(Myers et al., 2023)	Yes	22
(Metcalf et al., 2024)	Yes	50
(Myers et al., 2022)	Yes	50
(Christiano et al., 2017)	Yes	NA
(Hejna III & Sadigh, 2023)	Yes	4
(Hwang et al., 2023b)	Yes	5
(Biyik & Sadigh, 2018)	Yes	10
(Hwang et al., 2023a)	Yes	5
(Biyik et al., 2020)	Yes	10
(Wilde et al., 2021)	Yes	18
(Biyik et al., 2022)	Yes	15
(Ibarz et al., 2018)	Yes	NA
(Knox et al., 2023)	Yes	143
(Holk et al., 2024)	Yes	32
(Metcalf et al., 2023)	Yes	40
(Marta et al., 2024)	Yes	NA
(Marta et al., 2023a)	Yes	70
(Marta et al., 2023b)	Yes	20
(Ren et al., 2022)	Yes	NA
(Mehta & Losey, 2024)	Yes	15
(White et al., 2024)	Yes	20
(Wang et al., 2022)	Yes	10
(Wang et al., 2023)	No	
(Zhang & Kashima, 2024)	No	
(Lee et al., 2021a)	No	
(Liang et al., 2022)	No	
(Park et al., 2022)	No	
(Hu et al., 2024)	No	
(Liu et al., 2022)	No	
(Liu et al., 2023)	No	
(Xue et al., 2024)	No	
(Daniels-Koch & Freedman, 2022)	No	
(Verma et al., 2023)	No	
(Barnett et al., 2023)	No	
(Verma & Metcalf, 2024)	No	
(Metcalf et al., 2023)	No	
(Swamy et al., 2024)	No	
(Lee et al., 2021b)	No	
(Wang et al., 2021)	No	
(Maxence Hussonnois & Rana, 2023)	No	
(Wilson et al., 2012)	No	
(Liu et al., 2024)	No	
(Giovanelli et al., 2024)	No	
(Zhu et al., 2024)	No	
(Cheng et al., 2024)	No	

C ADDITIONAL HUMAN TEACHER RESULTS

Task	Feedback	Background	AUC	P Value	Final Performance	P Value
Cartpole-swingup	48	Non-CS	8258.44 ± 1898.47	0.183	648.88 ± 35.17	0.692
		CS	10208.11 ± 2311.16		614.14 ± 114.81	
Pendulum-swingup	40	Non-CS	7677.84 ± 3540.70	0.258	327.65 ± 407.80	0.23
		CS	10188.70 ± 3858.94		641.68 ± 237.76	

Table 7: This table shows the AUC and final performance of SDP (mean +/- 95% confidence intervals) for CS and non-CS participants in the human subject study.

Task	Feedback	Demographic	AUC	P Value
Cartpole-swingup	48	Female	8720.03 ± 3185.4	0.085
		Male	11770.22 ± 1669.96	
Pendulum-Swingup	40	Female	8587.13 ± 2772.96	0.331
		Male	10134.67 ± 4822.22	
Cartpole-swingup	48	Age 18-30	10109.54 ± 2370.61	0.313
		Age 50-70	10923.05 ± 1292.83	
Pendulum-Swingup	40	Age 18-30	8619.07 ± 3381.48	0.318
		Age 50-70	10607.93 ± 5249.72	
Cartpole-swingup	48	Education: Bachelors or higher	11044.55 ± 1783.25	0.794
		Education: High school diploma	6247.99 ± 5186.47	
Pendulum-Swingup	40	Education: Bachelors or higher	10188.88 ± 3858.95	0.742
		Education: High school diploma	7677.85 ± 3540.7	

Table 8: This table shows the AUC of SDP (mean +/- 95% confidence intervals) for different demographic conditions. This table highlights that SDP can be effective irrespective of demographic background (gender, age, and educational background).

Tables 7-8 highlight that there is no significant difference in the performance of SDP when human teachers have differing backgrounds. This includes a background in CS, educational level, age, and gender. We note that we did not perform statistical analysis comparing racial groups as some experiments had only one participant of a particular racial background, making comparisons less meaningful.

D ADDITIONAL SIMULATED TEACHER RESULTS

For simplicity, in all additional experiments in this section, we only compare SDP + PEBBLE with PEBBLE (or R-PEBBLE).

D.1 STUDY OF AGENT UPDATE PHASE

Figure 8 emphasizes how the agent update phase does result in new transitions, therefore the teacher provides feedback to transitions that are different from the original sub-optimal transitions used for pre-training.

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

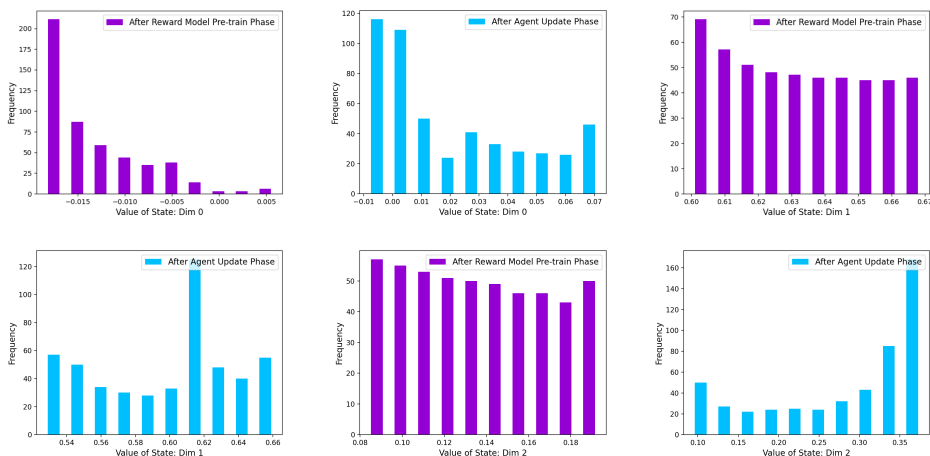


Figure 8: Door open, preference learning exp. These plots show how the agent’s policy has changed from the reward model pre-training phase (purple histograms) to the agent update phase (blue histograms), thereby resulting in a different distribution for the state features.

1312

1313

1314

1315

1316

1317

D.2 TRUE REWARD VALUES OF RANDOM POLICY

1318

1319

1320

1321

1322

1323

1324

D.3 ZERO WEIGHT STUDY

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

D.4 FAKE INPUT STUDY

1338

1339

1340

1341

1342

1343

1344

1345

1346

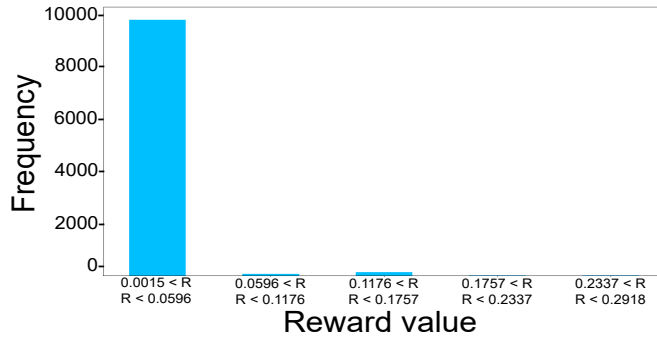
1347

1348

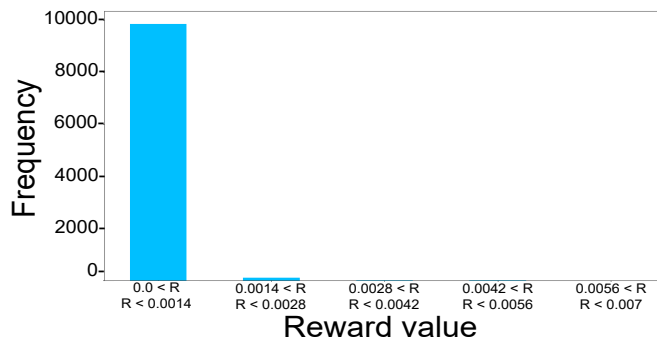
1349

Second, SDP makes use of state, action transitions that are gathered through a sub-optimal policy. Therefore, these are transitions that an agent experiences while interacting with its environment. However, as previously noted, the goal of the reward model pre-training phase is for the reward model to learn to output zero. Therefore, is it necessary for the reward model to pre-train on inputs that are real environment transitions? Instead, can we pre-train the reward model on transitions that did not result from an agent-environment interaction? To test this, we created “fake” inputs of size $\text{dim}(\text{state}) + \text{dim}(\text{action})$, and for each input dimension, we randomly sampled a value from $\mathcal{N}(0, 1)$. We obtained 50,000 “fake” transitions and used this data for the reward model pre-training phase. In this experiment, our goal is to understand the effect of the type of inputs on the reward model pre-training phase, therefore we kept the agent update phase as is (i.e., provided the true sub-optimal data for this phase). We then compared SDP using true transitions to SDP using “fake” transitions in Figures 10–left and 12. We found that the full SDP (purple curve) in which we pre-train the reward model on true sub-optimal transitions yields significantly greater final performance than SDP using

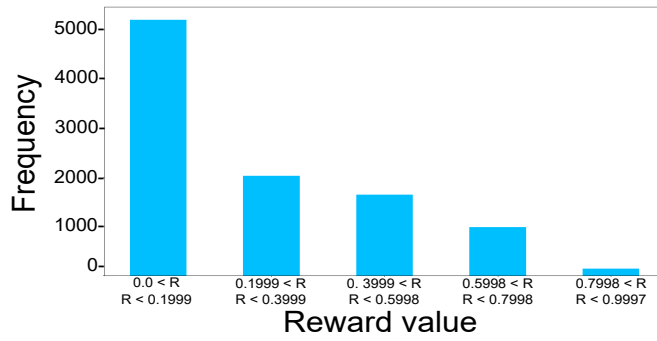
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403



(a) Walker-walk.



(b) Chetah-run.



(c) Quadruped-walk.

Figure 9: Distribution of true reward values for transitions obtained with a random policy

“fake” transitions (green curve). This highlights the importance of using true sub-optimal transitions in SDP.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

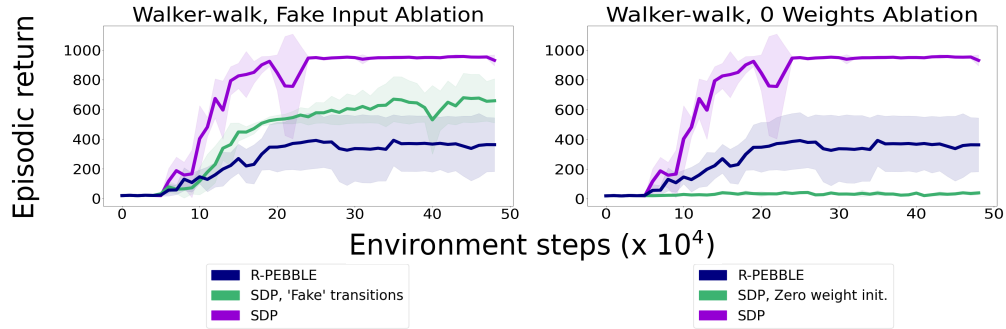
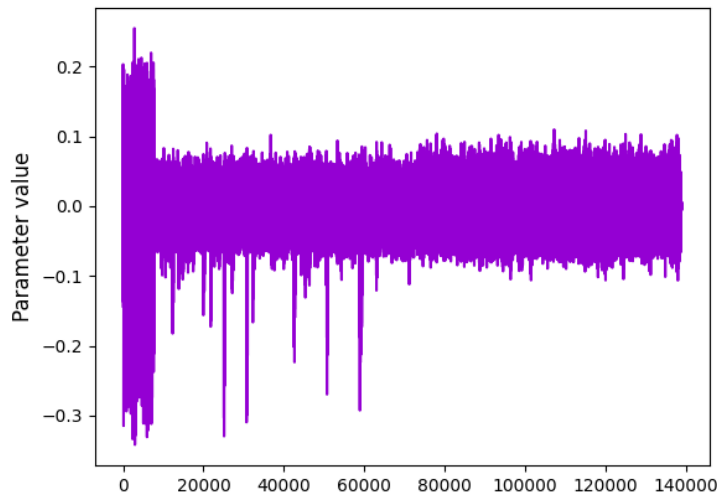


Figure 10: Zero weights and fake input studies in Walker-walk



(a) Scalar feedback experiment in Walker-walk.

Figure 11: This figure shows the reward model weights of SDP after the reward model pre-training phase. This demonstrates that the reward model pre-train phase of SDP does not result in zero neural network weights.

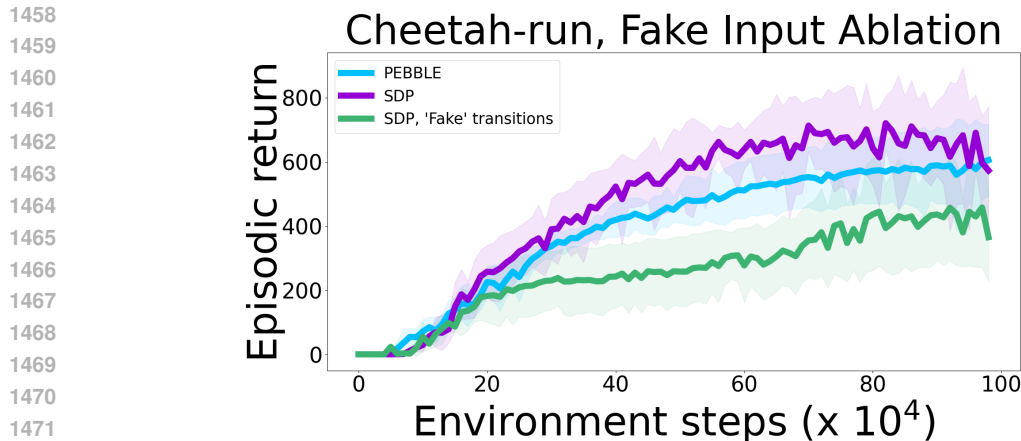


Figure 12: Fake Input Study in Cheetah Run Preference Learning

1473
1474
1475
1476

D.5 SDP COMPONENT STUDIES

1477
1478
1479
1480
1481
1482
1483
1484
1485

Figure 13 shows additional phase ablations in the preference learning experiments done on the Cheetah-run environment. This highlights the importance of using both phases in SDP. In addition, we ran an experiment in Walker-walk where we directly apply the procedure Yu et al. (2022) uses for leveraging sub-optimal data in the offline RL setting. They simply store the sub-optimal transitions in the RL agent’s replay buffer with the pseudo reward label of 0 and follow standard offline RL. We repeated this with the other difference of following standard preference learning (PEBBLE) afterwards. We found that the average final performance was 104.32 with a 95% confidence interval of 46.12, which is close to 4 times less than the final performance of SDP.

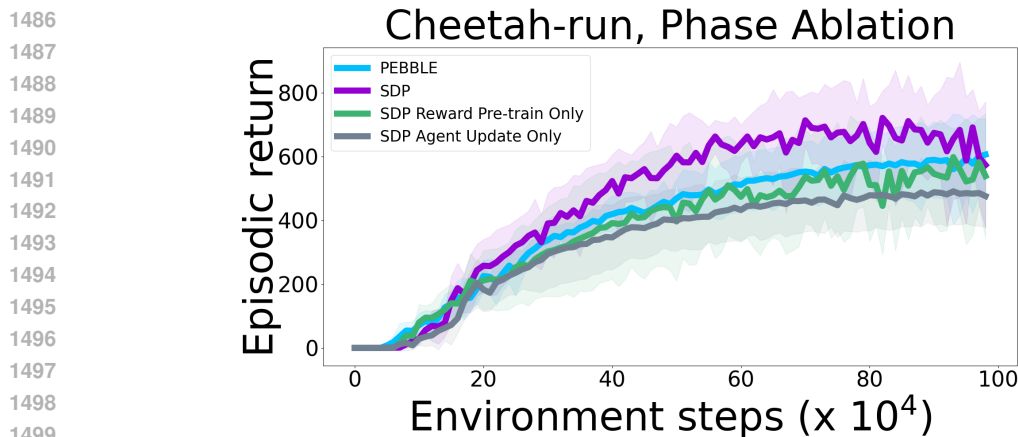


Figure 13: Phase Ablation Study in Cheetah Run Preference Learning

1500
1501
1502
1503
1504
1505

Figure 14 shows an additional ablation over the amount of prior data used in SDP. We observed similar performance gains as described in section 5.4.

1506
1507

D.6 EFFECT OF RELABELLING REPLAY BUFFER

1508
1509
1510
1511

SDP is combined with four preference learning algorithms, PEBBLE, RUNE, SURF, and MRN. A core feature of these algorithms is that every time the reward model is updated, all transitions inside the RL agent’s replay buffer are updated using the latest reward model. Figure 15 shows the effects of not relabelling the sub-optimal data (i.e., the reward labels) with the latest learned reward model. This means the reward labels remain frozen at zero throughout the entire training process.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

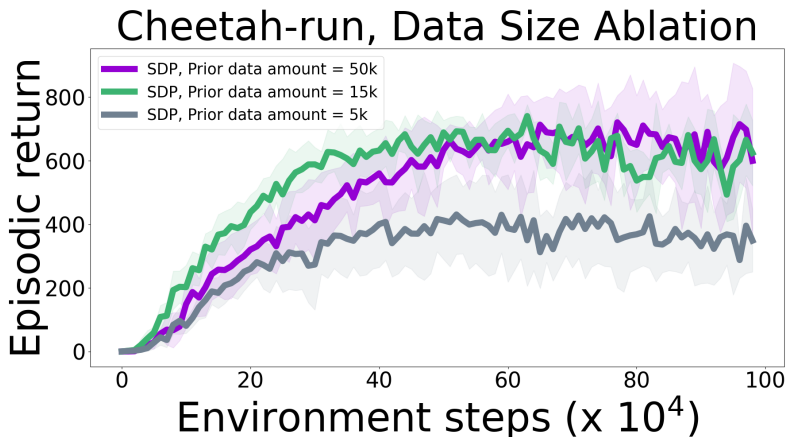


Figure 14: Prior Data Amount Study in Cheetah Run Preference Learning

This ablation was to demonstrate that if the sub-optimal data in the agent’s replay buffer is not updated with the latest reward model, then performance will suffer. This is likely the case because the incorrect reward bias from the pseudo-labeling process persists, whereas when we relabel the transitions with the updated reward model, the incorrect reward bias may reduce over time.

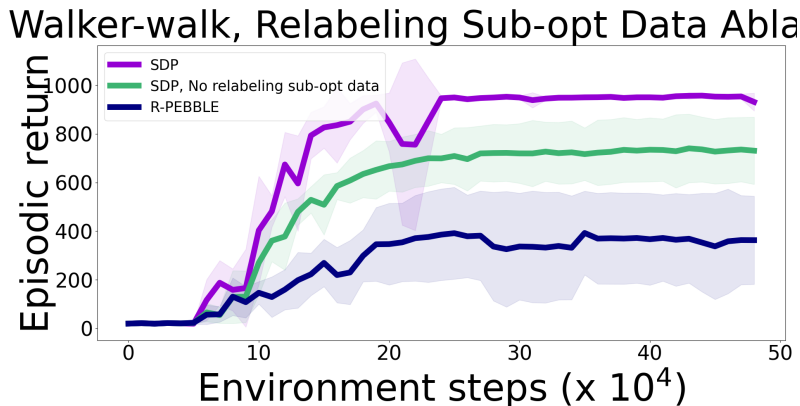


Figure 15: Relabeling sub-opt data study: Walker-walk, scalar feedback

D.7 EFFECT OF TRANSITION QUALITY IN SDP

Moreover, we show that the effectiveness of SDP relies on the use of sub-optimal data transitions. If we use high-quality data transitions (i.e., transitions that came from a fully trained RL agent policy), SDP will fail (see Figure 16). This unsurprising result confirms that pseudo-labeling high-reward transitions with zero can significantly hurt the reward model and the agent’s performance

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

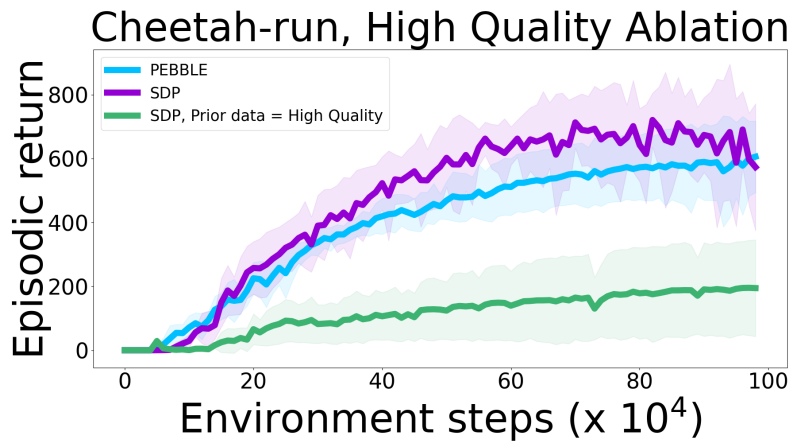


Figure 16: Using high-quality data in SDP study: Cheetah-run, preference feedback

1620 D.8 SCALAR-BASED EXPERIMENT STATISTICS
 1621

1622 Tables 9 and 10 provide a summary of the mean final performance and the mean area under the
 1623 curve for all environments and benchmarks in the scalar feedback setting
 1624

1625

Learning from scalar feedback			
	SDP	R-PEBBLE	Deep TAMER
1626 Walker-walk	931.31 \pm 36.59	362.99 \pm 181.23*	33.10 \pm 11.10*
1627 Cheetah-run	862.72 \pm 49.13	522.10 \pm 238.87*	30.71 \pm 16.26*
1628 Quadruped-walk	777.38 \pm 156.74	543.92 \pm 193.21*	73.74 \pm 48.38*

1629
1630

1631 Table 9: This table shows the mean final performance plus and minus the 95% confidence interval.
 1632 * indicates SDP achieves a significantly greater mean final performance.
 1633

1634

Learning from scalar feedback			
	SDP	R-PEBBLE	Deep TAMER
1635 Walker-walk	34981.65 \pm 1559.24	13147.1 \pm 6102.31*	1823.42 \pm 578.48*
1636 Cheetah-run	55144.89 \pm 4266.94	33557.58 \pm 15185.47*	2459.23 \pm 1226.34*
1637 Quadruped-walk	58370.85 \pm 10759.	37076.52 \pm 11968.48*	7365.68 \pm 3556.27*

1638
1639
1640

1641 Table 10: This table shows the mean area under the learning curve (AUC) plus and minus the 95%
 1642 confidence interval. * indicates SDP achieves a significantly greater AUC.
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

1674 D.9 PREFERENCE LEARNING EXPERIMENT STATISTICS
16751676 Tables 11-14 provide a summary of the mean final performance and the mean area under the curve
1677 for all environments and benchmarks in the preference feedback setting
1678

1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689
Task	Feedback	Method	Final Return	P Value						
1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691
			1692	1693	1694	1695	1696	1697	1698	1699
			1700	1701	1702	1703	1704	1705	1706	1707
			1708	1709	1710	1711	1712	1713	1714	1715
			1716	1717	1718	1719	1720	1721	1722	1723
			1724	1725	1726	1727				

1719 Table 11: This table shows the final performance (mean +/- 95% confidence intervals) for all DM-
1720 Control preference learning experiments. * indicates significant differences.
1721
1722
1723
1724
1725
1726
1727

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

Task	Feedback	Method	AUC	P Value
Cartpole-swingup	48	PEBBLE	8248.14 ± 2148.59	0.0*
		SDP + PEBBLE	21519.9 ± 2435.02	
		RUNE	8817.21 ± 3534.57	0.004*
		SDP + RUNE	20970.54 ± 4643.96	
		SURF	5593.44 ± 1099.56	0.074
		SDP + SURF	10288.22 ± 4584.16	
		MRN	6774.95 ± 1783.95	
		SDP + MRN	22564.32 ± 3846.6	0.0*
Walker-walk	100	PEBBLE	7571.79 ± 2903.09	0.001*
		SDP + PEBBLE	18458.81 ± 1224.69	
		RUNE	6525.52 ± 2566.88	0.0*
		SDP + RUNE	21059.76 ± 3056.28	
		SURF	7675.24 ± 2857.99	0.001*
		SDP + SURF	16700.51 ± 1449.62	
		MRN	7711.05 ± 1295.46	
		SDP + MRN	25184.26 ± 3071.65	0.0*
Cheetah-run	200	PEBBLE	39245.94 ± 5189.16	0.006*
		SDP + PEBBLE	51121.68 ± 3488.87	
		RUNE	45824.63 ± 4459.04	0.792
		SDP + RUNE	37202.31 ± 16330.83	
		SURF	40522.18 ± 2989.22	0.218
		SDP + SURF	42777.1 ± 3775.07	
		MRN	47569.45 ± 2191.57	
		SDP + MRN	49968.24 ± 2187.97	0.106
Quadruped-walk	500	PEBBLE	20584.95 ± 9814.49	0.018*
		SDP + PEBBLE	42044.34 ± 11183.84	
		RUNE	22861.89 ± 9861.31	0.082
		SDP + RUNE	32664.79 ± 3929.32	
		SURF	30884.8 ± 11660.52	0.074
		SDP + SURF	44176.44 ± 8404.76	
		MRN	14805.07 ± 1860.39	
		SDP + MRN	35209.32 ± 6225.09	0.002*

Table 12: This table shows the AUC (mean +/- 95% confidence intervals) for all DMControl preference learning experiments. * indicates significant differences.

Task	Feedback	Method	Final Return	P Value
Hammer	7500	PEBBLE	10.0 ± 13.58	0.003*
		SDP + PEBBLE	66.0 ± 21.18	
		RUNE	4.0 ± 7.01	0.001*
		SDP + RUNE	68.0 ± 17.0	
		SURF	0.0 ± 0.0	0.033*
		SDP + SURF	36.0 ± 25.16	
		MRN	4.0 ± 7.01	0.027*
		SDP + MRN	46.0 ± 27.5	
Door-unlock	500	PEBBLE	42.0 ± 34.35	0.066
		SDP + PEBBLE	80.0 ± 15.68	
		RUNE	26.0 ± 22.59	0.354
		SDP + RUNE	36.0 ± 38.65	
		SURF	20.0 ± 35.06	0.02*
		SDP + SURF	80.0 ± 22.17	
		MRN	48.0 ± 38.17	0.399
		SDP + MRN	56.0 ± 37.02	
Door-lock	500	PEBBLE	38.0 ± 30.06	0.035*
		SDP + PEBBLE	80.0 ± 7.84	
		RUNE	62.0 ± 32.51	0.793
		SDP + RUNE	40.0 ± 30.87	
		SURF	70.0 ± 31.85	0.684
		SDP + SURF	60.0 ± 13.58	
		MRN	66.0 ± 26.93	0.268
		SDP + MRN	78.0 ± 17.88	
Drawer-open	1000	PEBBLE	4.0 ± 7.01	0.045*
		SDP + PEBBLE	36.0 ± 25.16	
		RUNE	0.0 ± 0.0	0.001*
		SDP + RUNE	66.0 ± 14.24	
		SURF	20.0 ± 35.06	0.018*
		SDP + SURF	80.0 ± 14.67	
		MRN	20.0 ± 35.06	0.136
		SDP + MRN	56.0 ± 40.21	
Window-open	200	PEBBLE	34.0 ± 37.44	0.234
		SDP + PEBBLE	54.0 ± 26.35	
		RUNE	12.0 ± 21.04	0.098
		SDP + RUNE	38.0 ± 24.42	
		SURF	14.0 ± 20.44	0.014*
		SDP + SURF	66.0 ± 26.35	
		MRN	46.0 ± 32.61	0.208
		SDP + MRN	68.0 ± 31.06	

Table 13: This table shows the final performance (mean +/- 95% confidence intervals) for all Meta-world preference learning experiments. * indicates significant differences.

Task	Feedback	Method	AUC	P Value
Hammer	7500	PEBBLE	424.0 ± 175.18	0.017*
		SDP + PEBBLE	2222.0 ± 1009.47	
		RUNE	490.0 ± 300.85	0.001*
		SDP + RUNE	2520.0 ± 627.47	
		SURF	482.0 ± 215.59	0.006*
		SDP + SURF	3042.0 ± 1073.54	
		MRN	494.0 ± 311.76	0.016*
		SDP + MRN	2494.0 ± 1108.52	
Door-unlock	500	PEBBLE	2142.0 ± 1738.06	0.031*
		SDP + PEBBLE	4838.0 ± 1251.13	
		RUNE	1272.0 ± 1275.71	0.348
		SDP + RUNE	1836.0 ± 2063.68	
		SURF	1580.0 ± 2462.83	0.047*
		SDP + SURF	4802.0 ± 1537.81	
		MRN	3082.0 ± 2232.43	0.342
		SDP + MRN	3842.0 ± 2230.38	
Door-lock	500	PEBBLE	876.0 ± 647.87	0.008*
		SDP + PEBBLE	2300.0 ± 135.0	
		RUNE	1342.0 ± 590.05	0.809
		SDP + RUNE	884.0 ± 633.87	
		SURF	1892.0 ± 842.85	0.429
		SDP + SURF	1986.0 ± 184.83	
		MRN	1966.0 ± 751.51	0.306
		SDP + MRN	2236.0 ± 479.9	
Drawer-open	1000	PEBBLE	314.0 ± 533.02	0.117
		SDP + PEBBLE	1702.0 ± 1706.12	
		RUNE	38.0 ± 49.71	0.004*
		SDP + RUNE	3794.0 ± 1306.6	
		SURF	696.0 ± 1167.88	0.015*
		SDP + SURF	4078.0 ± 1839.07	
		MRN	1550.0 ± 2690.98	0.211
		SDP + MRN	3206.0 ± 2114.16	
Window-open	200	PEBBLE	794.0 ± 1015.04	0.131
		SDP + PEBBLE	1698.0 ± 822.8	
		RUNE	230.0 ± 334.14	0.064
		SDP + RUNE	948.0 ± 631.92	
		SURF	242.0 ± 248.44	0.01*
		SDP + SURF	2048.0 ± 867.93	
		MRN	1252.0 ± 796.8	0.223
		SDP + MRN	1728.0 ± 666.48	

Table 14: This table shows the area under the curve (mean +/- 95% confidence intervals) for all Metaworld preference learning experiments. * indicates significant differences.