

---

# Exploring and Addressing Reward Confusion in Offline Preference Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Spurious correlations in a reward model’s training data can prevent Reinforcement  
2 Learning from Human Feedback (RLHF) from identifying the desired goal and  
3 induce unwanted behaviors. In this work, we study the reward confusion problem  
4 in offline RLHF where spurious correlations exist in data. We create a lightweight  
5 benchmark to study this problem and propose a method that can reduce reward  
6 confusion by leveraging model uncertainty and the transitivity of preferences with  
7 active learning.

## 1 Introduction

9 For many real-world tasks, designing adequate reward functions is challenging, which has led to the  
10 rise of Reinforcement Learning from Human Feedback (RLHF) [6]. In this work, we study a failure  
11 mode of offline RLHF that we refer to as *reward confusion*. This occurs when the reward  $R$  in a Markov  
12 Decision Process (MDP) is a function of features  
13  $z_1, \dots, z_n$  inferred from the observation-action pair  
14  $(o, a)$ . In a simplified scenario,  $R$  relies on  $z_1$  but not  
15  $z_2$ , yet  $z_1$  and  $z_2$  are highly correlated in the training  
16 data. An empirical risk minimizer might mistakenly  
17 conclude that  $z_2$  affects  $R$ . As we’ll see, this incorrect  
18 dependence can lead to failures when training  
19 a policy against the learned reward function. We  
20 graphically illustrate this problem in Figure 1.

24 To better understand this phenomenon, we created a benchmark environment called *Confusing*  
25 *Minigrid (CMG)* that tests reward confusion in models. We carefully designed six tasks with three  
26 types of spurious information for the minigrid environment, which we introduce in detail in Appendix  
27 A. We will open source the benchmark’s code soon.

28 Besides the CMG benchmark, one other our major contributions is an algorithm named Information-  
29 Guided Preference Chain (IMPEC) designed to address the reward confusion problem. It involves  
30 two stages of training: First, we use information gain as the acquisition function to select comparison  
31 rollouts that reduce uncertainty about the reward function. Second, we form a complete preference  
32 ordering over the set of selected rollouts, rather than just a partial ordering as in traditional RLHF.

33 Our experiments show that these techniques together improve sample efficiency while reducing  
34 reward confusion. We show in Section 4 that using the same comparison budget, IMPEC can

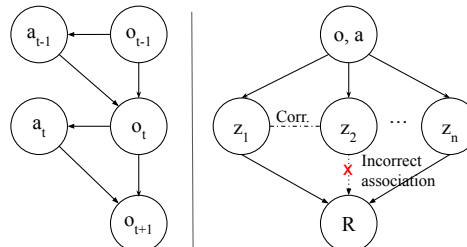


Figure 1: Illustration of a simplified MDP (left) and reward confusion (right). The reward  $R$  is a function of the feature  $z_1$ , but not  $z_2$ . Spurious correlation between  $z_1$  and  $z_2$  can cause a network to wrongly model  $R$  as a function of  $z_2$ .

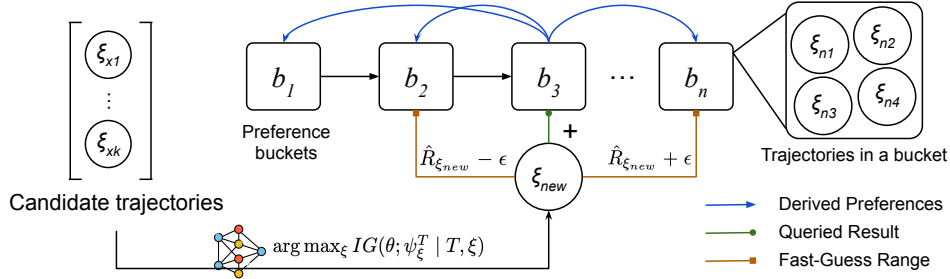


Figure 2: The IMPEC algorithm creates a sorted preference chain of  $n$  buckets, each containing one or more rollouts with equal returns.

35 outperform many other active preference learning baselines. To the best of our knowledge, it is the  
 36 first algorithm that attempts to solve the reward confusion problem in preference learning.

## 37 2 Related Work

38 **Causal Confusion** The problem of *causal confusion*, which refers to models learning to depend on  
 39 spurious correlations in the training data, has been studied in behavioral cloning [7], reinforcement  
 40 learning [12], and reward learning [20]. Past work shows empirically and theoretically that spurious  
 41 correlations and confounders in the training set can worsen an agent’s deployment performance [22, 9].  
 42 Reward confusion is essentially causal confusion that occurs during reward learning.

43 **Goal Misgeneralization** While past work on causal confusion studies it as a cause of complete  
 44 failure to learn goal-directed behavior, it can also make agents optimize for incorrect goals, i.e. goal  
 45 misgeneralization. For example, in Procgen’s CoinRun, the coin to be picked up is always on the  
 46 right. RL agents can confuse “running to the right” with the real goal of “getting the coin” [10].  
 47 Generally, a model’s behavior can be consistent with a goal, but it may not be the test-time goal [16].

48 **Preference Learning** Learning reward models from preference labels [6] have gained traction due  
 49 to their low cost compared to expert demonstrations [23] or language inputs [21]. We have also seen  
 50 progress in other tasks of “reward engineering”, e.g. reward hacking [18].

## 51 3 Method

52 **Models** We consider an agent in an environment following the Markov Decision Process (MDP)  
 53 defined by  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ .  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, \infty)$  is  
 54 the transition probability density. A rollout  $\xi = (s_t, a_t)$  is a sequence of states and actions. Given  
 55 unranked rollouts  $\Xi$ , our algorithm actively collects ranking information to sort them into an ordered  
 56 list  $T = \langle \xi_1, \xi_2, \dots, \xi_n \rangle$ . The rank of rollout  $\xi \in T$  is denoted by  $\psi_{\xi}^T$ . On each transition, the  
 57 environment emits a reward  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$ . Our goal is to obtain  $\mathcal{R}^*$  that induces correct policies.

58 **Preferences** We model the human’s probability of preferring  $\xi_1$  in a pair  $(\xi_1, \xi_2)$  through the  
 59 Shepard-Luce choice rule [17, 13]:  $P[\xi_1 \succ \xi_2] = \frac{\exp \sum_t r(o_t^1, a_t^1)}{\exp \sum_t r(o_t^1, a_t^1) + \exp \sum_t r(o_t^2, a_t^2)}$ . We extend this  
 60 model to a ternary one by allowing the human to flag when two rollouts are equally good,  $\xi_1 \equiv \xi_2$ .  
 61 We use cross-entropy loss to improve reward model’s predictions for human’s true preference.

### 62 3.1 Information-Guided Preference Chain (IMPEC)

63 **Key intuition: Increase Contrast Among Valuable Rollouts.** In most preference comparison  
 64 algorithms, a rollout  $\xi_1$ ’s relation is considered explicitly only with another one  $\xi_2$ . Suppose that in  
 65 the ground truth,  $\xi_1 \succ \xi_2 \succ \xi_3 \succ \xi_4$ , and we already know  $\xi_1 \succ \xi_2, \xi_3 \succ \xi_4$ . To figure out  $\xi_1$  and  
 66  $\xi_2$ ’s relationship with  $\xi_3$  and  $\xi_4$ , the most efficient query is whether  $\xi_2 \succ \xi_3$ . Once we establish that,  
 67 we can immediately obtain the preference relations on all four rollouts.

68 **Creating and Maintaining a Preference Chain** We maintain an ordered chain for rollouts. Starting  
 69 from an empty chain, for each new rollout we queried from the dataset, we imitate insertion sort by  
 70 recursively finding the ranking of it using human’s preference labels. Hence, by the time we observe

71 all the rollouts, we have a sorted list of rollouts, ordered according to human preferences. Rollouts  
 72 can have identical returns, so we treat each element of the chain as a *bucket*  $b \in \mathcal{B}$  of rollouts with the  
 73 same return. If the human decides that a new rollout  $\xi_{\text{new}}$  is equally preferred to  $\xi_m$  in bucket  $b_m$ ,  
 74 then  $\xi_{\text{new}}$  will be added to  $b_m$ . On the other hand, if  $\xi_{\text{new}} \succ \xi_m$  and  $\xi_{\text{new}} \prec \xi_{m-1}$  ( $\xi_{m-1}$  resides in a  
 75 previous bucket  $b_{m-1}$ ), then the algorithm will insert a new bucket containing only  $\xi_{\text{new}}$  in between  
 76  $b_m$  and  $b_{m-1}$ . This ensures that  $b_0$  contains the best rollouts seen so far and  $b_n$  contains the least  
 77 preferred rollouts (where  $n$  is the chain length). We illustrate this process in Figure 2.

78 Our reward model is a Bayesian neural network  
 79 (BNN) [3] which maintains a Gaussian distribu-  
 80 tion over a network’s weights and biases. As we  
 81 will see, this allows us to incorporate epistemic  
 82 uncertainty over reward functions into the ac-  
 83 tive selection procedure. In the noiseless case,  
 84 insertion sort needs  $O(\log n)$  queries to find the  
 85 position for  $\xi_{\text{new}}$ . However, we have access to  
 86 a partially trained reward network, which we  
 87 use to guess the rank for  $\xi_{\text{new}}$ , reducing the  
 88 number of buckets we must search over. We  
 89 include more design details in Appendix B for  
 90 the design of the fast query.

91 **Information Gain** Given an existing chain  
 92 of rollouts, we use information gain as the ac-  
 93 quisition function to decide which rollouts to  
 94 compare next, so we reduce the most uncertainty over network weights. The information gain over  
 95 network weights  $\theta$  by selecting a rollout  $\xi \in \Xi$  for ranking is

$$I(\theta; \psi_\xi^T | T, \xi) = H(\theta | T, \xi) - H(\theta | \psi_\xi^T, T, \xi) \quad (1)$$

96 where  $\psi_\xi^T$  is the rollout’s ranking on chain  $T$ . Intuitively, it measures how much we expect to  
 97 reduce uncertainty about the weights after observing the ranking  $\psi_\xi^T$  of rollout  $\xi$ . As shown in  
 98 Appendix C, Equation 1 (information gain) can be approximated by drawing  $M$  weight samples,  
 99  $\theta_1, \theta_2, \dots, \theta_M \sim \theta$ , from the posterior through

$$\frac{1}{M} \sum_{i=1}^M \sum_{\psi} P(\psi_\xi^T | T, \theta_i, \xi) \cdot \log \left( \frac{M \cdot P(\psi_\xi^T | T, \theta_i, \xi)}{\sum_{\theta_j} P(\psi_\xi^T | T, \theta_j, \xi)} \right) \quad (2)$$

100  $P(\psi_\xi^T | T, \theta, \xi)$  is a complicated distribution, and so we (loosely) approximate it with Equation 3.  
 101 Intuitively, it is proportional to the probability that  $\xi_i \succ \xi \succ \xi_{i+1}$ .

$$P(\psi_\xi^T = i | T, \theta, \xi) \propto P(\xi_i \succ \xi | \theta) \cdot P(\xi \succ \xi_{i+1} | \theta) \quad (3)$$

102 We summarize the complete process in Algorithm 1. The network is first supervised trained on the  
 103 preference dataset  $D$  using cross entropy loss with  $P[\xi_1 \succ \xi_2]$  modeled through the Shepard-Luce  
 104 choice rule. With the limited query budget for human preferences, we first find out the rollout  $\xi$   
 105 whose ranking  $\psi_\xi^T$  on chain  $T$  will provide the most information gain over the model weights  $\theta$ . Then  
 106 we use insertion sort to find out  $\xi$ ’s real ranking  $\psi_\xi^{T*}$  in the chain. We add the rollout  $\xi$  onto the  
 107 appropriate position of the chain  $T$ , then based on its position, derive preference labels with all other  
 108 rollouts on the chain. We repeat this process until the network weight has converged.

## 109 4 Experiments

110 **Experiment Settings** We compare our method with the standard RLHF algorithm, and two other  
 111 RLHF with active learning methods: pairwise information gain [2] and pairwise volume removal [14].  
 112 The information gain (IG) method is similar to ours but reasons only about individual preference  
 113 pairs and not about the result of the ranking process. The volume removal method was designed in

---

### Algorithm 1 The IMPEC Algorithm

---

**Require:** Preference dataset  $D$ , network  $\theta$ , query  
 budget  $Q$   
 $T \leftarrow []$   
**while** not converged **do**  
 $\theta \leftarrow \text{SupervisedTrain}(\theta, D)$   
**if** budget not reached **then**  
 $\xi \leftarrow \arg \max_{\xi} I(\theta; \psi_\xi^T | T, \xi)$   
 $\psi_\xi^{T*} \leftarrow \text{InsertionSort}(\xi, T, \theta)$   
 $T \leftarrow T \cup \xi$   
 $D \leftarrow D \cup \text{DerivePreferences}(\xi, T, \psi_\xi^T)$   
**end if**  
 $i \leftarrow i + 1$   
**end while**

---

	Baseline	IMPEC	Infogain	Vol Removal
Empty	12.90±8.20	18.13±2.15 <sup>†</sup>	7.33±9.98	14.66±8.41
Dyn Obs	7.17±6.56	11.78±2.82 <sup>†</sup>	5.34±6.75	4.06±6.03
Lava	8.70±12.83	17.65±1.97 <sup>†</sup>	13.33±8.35	9.75±9.39
Lava Pos.	3.51±1.26	7.21±2.20 <sup>†</sup>	3.60±1.52	4.66±2.53
Fetch	10.60±2.85	11.52±1.17	9.93±2.47	10.80±1.51
Door	1.58±0.31	1.67±0.53	1.55±0.48	1.68±0.71

Table 1: Ground truth returns (mean  $\pm$  standard deviation) for different methods on Confusing Minigrid. <sup>†</sup> indicates a p-value of  $\leq 0.1$  (vs. baseline).

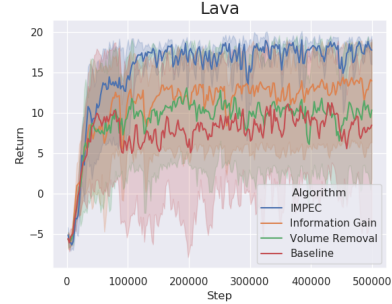


Figure 3: The learning curves for 4 algorithms trained on the Lava task.

114 the linear reward setting to reduce the volume of weight vectors supported under the posterior after  
 115 each preference update. We conduct experiments on 6 CMG tasks. Detailed information on each task  
 116 and their added spurious correlations can be found in Appendix A.

117 We first perform offline reward learning, and then apply online reinforcement learning using the  
 118 learned reward function to obtain a policy. The RL agent receives rewards from the learned function  
 119 instead of the environment, and is trained with Proximal Policy Optimization [15]. More detailed  
 120 experiment and hyperparameter settings can be found in Appendices D and E.

121 **Main Results** Performance on all six tasks can be found in Table 1, with each run repeated over  
 122 five seeds. We further compute the p-values of the results being better than baseline performance,  
 123 with complete results in Appendix H. Except for the task Go To Door where all algorithms perform  
 124 poorly, IMPEC has a higher mean return than other algorithms, and often a lower standard deviation.  
 125 Although IMPEC achieves a higher mean than the other methods on most tasks, its p-values are  $\leq 0.1$   
 126 (per appendix H). This provides some (albeit not particularly strong) statistical evidence that IMPEC  
 127 has a better mean performance distribution from the baseline. In contrast, the volume removal and  
 128 information gain methods are often statistically indistinguishable from the baseline.

129 The decisive factor for each algorithm’s performance is how often they fail: Out of the 5 seeds, how  
 130 often does the reward function learn to optimize for the spurious goal? All methods but IMPEC have  
 131 fairly high probability of taking the spurious feature as the “correct” feature, and hence rewarding  
 132 incorrect behaviors and obtaining low ground truth returns. We plot the learning curves for each  
 133 algorithm in the Lava task, where the baseline curve has the largest standard deviations. We observe  
 134 similar phenomena in other tasks. We include an ablation study in Appendix F. The complete learning  
 135 curves are included in Appendix J.

136 **Limitations and Future Work** A limitation of IMPEC is its potential sensitivity to noise in  
 137 preferences. In our experiments, we keep a relatively low noise level, and we believe more sophisti-  
 138 cated algorithms could improve robustness to noise, perhaps inspired by past work on noisy binary  
 139 search [8, 5]. Our results suggest a deeper connection between the quality of preference datasets  
 140 and the efficiency of preference learning algorithms. In appendix G, we show some first steps of a  
 141 graph theoretic analysis for reasoning about preference dataset quality. We are interested in further  
 142 exploring the influence of graph-theoretic qualities and their effects on preference learning, and using  
 143 the insights in future algorithm design.

## 144 5 Conclusion

145 This work studies the reward confusion problem. Our experiments on the Confusing Minigrid  
 146 benchmark show that reward confusion in offline preference learning can lead to undesired policy  
 147 behaviors. The benchmark is easy to configure, and we expect it to be particularly useful for iterative  
 148 research. In addition, we proposed IMPEC to reduce the impact of reward confusion. It exploits  
 149 preference transitivity and obtains decent empirical performance on tasks with different sources of  
 150 reward confusion. We believe that the findings of our work will be helpful for making AI more  
 151 aligned with human values.

## References

- 152
- 153 [1] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning  
154 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*  
155 *Research*, 47:253–279, 2013.
- 156 [2] Erdem Bıyık, Malayandi Palan, Nicholas C Landolfi, Dylan P Losey, and Dorsa Sadigh.  
157 Asking easy questions: A user-friendly approach to active reward learning. *arXiv preprint*  
158 *arXiv:1910.04365*, 2019.
- 159 [3] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty  
160 in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR,  
161 2015.
- 162 [4] Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld environ-  
163 nment for gymnasium, 2018.
- 164 [5] Sung-En Chiu. *Noisy binary search: Practical algorithms and applications*. University of  
165 California, San Diego, 2019.
- 166 [6] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
167 reinforcement learning from human preferences. *Advances in neural information processing*  
168 *systems*, 30, 2017.
- 169 [7] Pim De Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning.  
170 *Advances in Neural Information Processing Systems*, 32, 2019.
- 171 [8] Richard M Karp and Robert Kleinberg. Noisy binary search and its applications. In *Proceed-*  
172 *ings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 881–890.  
173 Citeseer, 2007.
- 174 [9] Daniel Kumor, Junzhe Zhang, and Elias Bareinboim. Sequential causal imitation learning with  
175 unobserved confounders. *Advances in Neural Information Processing Systems*, 34:14669–14680,  
176 2021.
- 177 [10] Lauro Langosco, Jack Koch, Lee D Sharkey, Jacob Pfau, and David Krueger. Goal misgener-  
178 alization in deep reinforcement learning. In *International Conference on Machine Learning*,  
179 pages 12004–12019. PMLR, 2022.
- 180 [11] Vito Latora and Massimo Marchiori. Efficient behavior of small-world networks. *Physical*  
181 *review letters*, 87(19):198701, 2001.
- 182 [12] Minne Li, Mengyue Yang, Furui Liu, Xu Chen, Zhitang Chen, and Jun Wang. Causal world  
183 models by unsupervised deconfounding of physical dynamics. *arXiv preprint arXiv:2012.14228*,  
184 2020.
- 185 [13] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. John Wiley & Sons, Inc.,  
186 1959.
- 187 [14] Dorsa Sadigh, Anca D Dragan, Shankar Sastry, and Sanjit A Seshia. *Active preference-based*  
188 *learning of reward functions*. 2017.
- 189 [15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
190 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 191 [16] Rohin Shah, Vikrant Varma, Ramana Kumar, Mary Phuong, Victoria Krakovna, Jonathan  
192 Uesato, and Zac Kenton. Goal misgeneralization: Why correct specifications aren’t enough for  
193 correct goals. *arXiv preprint arXiv:2210.01790*, 2022.
- 194 [17] Roger N Shepard. Stimulus and response generalization: A stochastic model relating general-  
195 ization to distance in psychological space. *Psychometrika*, 22(4):325–345, 1957.
- 196 [18] Joar Skalse, Nikolaus HR Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and  
197 characterizing reward hacking. *arXiv preprint arXiv:2209.13085*, 2022.

- 198 [19] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David  
 199 Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite.  
 200 *arXiv preprint arXiv:1801.00690*, 2018.
- 201 [20] Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca D Dragan, and Daniel Brown. A study  
 202 of causal confusion in preference-based reward learning. *arXiv preprint arXiv:2204.06601*,  
 203 2022.
- 204 [21] Hsiao-Yu Tung, Adam W Harley, Liang-Kang Huang, and Katerina Fragkiadaki. Reward  
 205 learning from narrated demonstrations. In *Proceedings of the IEEE Conference on Computer  
 206 Vision and Pattern Recognition*, pages 7004–7013, 2018.
- 207 [22] Junzhe Zhang, Daniel Kumor, and Elias Bareinboim. Causal imitation learning with unobserved  
 208 confounders. *Advances in neural information processing systems*, 33:12263–12274, 2020.
- 209 [23] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy  
 210 inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.

## 211 A Confusing Minigrid Definition

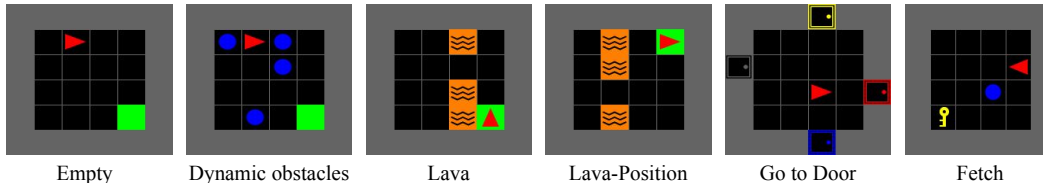


Figure 4: The six Confusing Minigrid tasks. The tasks require agents to move to goal positions, go to doors, or fetch objects.

212 Reward learning methods are often evaluated on benchmarks originally developed for reinforcement  
 213 learning algorithms, like Atari [1] or MuJoCo continuous control [19]. These reinforcement learning  
 214 benchmarks cannot easily be used to test reward confusion failures. Thus we created six new tasks  
 215 based on Minigrid [4], which together form the Confusing Minigrid benchmark.

Task name	Extra	Position	Color	DistShift
Empty	Y	N*	N*	N
Dynamic Obs.	Y	N*	N*	N
Lava	Y	N*	N*	N
Lava-Position	N*	Y	N*	Y
Go to door	N*	Y	N*	Y
Fetch	N*	N*	Y	Y

Table 2: A summary of Confusing Minigrid tasks. Y/N is short for yes/no. An asterisk means that a feature is off by default but can be optionally turned on. The “extra” column refers for having extra observation dimensions for the task. The position and color columns refer to spurious correlations that come from position and color information, respectively. DistShift means there is a difference between the training and testing environments.

### 216 A.1 Spurious Correlations from Extra Observations

217 This set of tasks tests whether spurious correlations with redundant observation dimensions can  
 218 interfere with learning. In these tasks, an agent needs to navigate to the goal cell. It can observe  
 219 the state of a glass of water it is holding, which exhibits “ripples” as it moves, and calms down if  
 220 the agent’s position remains unchanged. On rollouts where the agent moves straight to the goal and  
 221 stops, the level of ripples will be predictive of whether the agent has reached the goal, even though in  
 222 general it is possible to cause the ripples to disappear by stopping in any location and not just at the  
 223 goal. The training and testing variants of these tasks are the same.

224 **Empty** This is the simplest task in the benchmark. The agent needs to navigate to the goal cell, with  
225 all other cells being empty. The environment gives a positive reward when the agent reaches the goal,  
226 and zero otherwise.

227 **Dynamic Obstacles** This task augments the Empty task with obstacles that show up randomly in  
228 non-goal cells at each time step, which block the agent’s path. We include this task because the  
229 obstacles may stop the agent at a random grid while it is collecting rollouts, which increases the  
230 chance that the reward learning algorithm later realizes that “stabilizing the water” is not the correct  
231 goal.

232 **Lava** The environment contains a row of lava cells with a gap that allows the agent to cross and  
233 reach the goal. Standing on the lava cells gives a reward of -1. We use this task to observe if negative  
234 rewards will have any impacts on learning spurious correlations, and if the correlation is exploited,  
235 whether the agent will at least avoid the lava (low reward) area.

## 236 A.2 Spurious Correlations from Distributional Shifts

237 We design these tasks with different training and testing variants. The training environments have a  
238 90% probability where the goal configuration is spurious.

239 **Lava-Position** It is a variant of Lava with changing goal positions. In the training variant, the goal  
240 cell is usually located at one particular location. In the test variant, the goal grid can appear at other  
241 locations too.

242 **Go to Door** In this task, an agent is asked to move to a position adjacent to the goal door embedded  
243 in one of the four walls surrounding the grid. There are always four doors in the environment, and  
244 the goal door is most likely to be placed in the upper wall during training. During testing, the goal  
245 door can be placed in any of the four walls.

246 **Fetch** The agent’s goal in this task is to pick up a key. There is usually a distractor object in the  
247 environment that an agent can also pick up. In the training variant of this task, most keys are yellow,  
248 and most distractor objects are non-yellow. At the test time, the keys and distractor objects can appear  
249 in any color with equal chance.

250 For all these tasks, the agent receives a reward of +1 when the goal condition is satisfied. We  
251 summarize the task settings in Table 2. Changing the confounding type in Confusing Minigrad is  
252 as easy as modifying a keyword argument when initializing the environment. Training a standard  
253 preference learning algorithm on the tasks with the simplest observation type takes around 2.5 hours  
254 on a single Nvidia A6000 GPU, which is faster than many other image-based environments or  
255 complex control tasks.

## 256 B Use a Partially Trained Model to Reduce Queries

257 The network will first update its reward predictions on rollouts in each bucket  $\hat{R}_{b_i}$ , then predict  
258 the incoming rollout’s reward  $\hat{R}_{\xi_{\text{new}}}$ . Instead of using a point estimate  $\hat{R}_{\xi_{\text{new}}}$ , we use an interval  
259  $[\hat{R}_{\xi_{\text{new}}} - \epsilon, \hat{R}_{\xi_{\text{new}}} + \epsilon]$  to fast-guess where  $\xi_{\text{new}}$  may belong. IMPEC queries human preferences of  
260 two pairs  $(\xi_{\text{new}}, \xi_l)$  and  $(\xi_{\text{new}}, \xi_u)$ .  $\xi_l$  and  $\xi_u$  are the rollouts that are closest to the prediction lower  
261 bound  $\hat{R}_{\xi_{\text{new}}} - \epsilon$  and the upper bound  $\hat{R}_{\xi_{\text{new}}} + \epsilon$ , respectively. The  $\epsilon$  we use is the standard deviation  
262 of  $\hat{R}_{\xi_{\text{new}}}$  by  $M$  samples of the BNN. After  $\xi_{\text{new}}$  is added into the chain, we can derive preference  
263 relations between it and other rollouts by transitivity.

264 **C The Information Gain Objective Derivation**

265 This derivation is adapted from [2].

$$\begin{aligned}
I(\theta; \psi_\xi^T \mid T, \xi) &= H(\theta \mid T, \xi) - H(\theta \mid \psi_\xi^T, T, \xi) \\
&= -\mathbb{E}_{\theta, T, \xi} [\log P(\theta \mid T, \xi)] + \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} [\log P(\theta \mid \psi_\xi^T, T, \xi)] \\
&= -\mathbb{E}_{\theta, \psi_\xi^T, \xi, T} [\log P(\theta \mid T, \xi)] + \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} [\log P(\theta \mid \psi_\xi^T, T, \xi)] \\
&= \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} [\log P(\theta \mid \psi_\xi^T, T, \xi) - \log P(\theta \mid T, \xi)] \\
&= \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} \left[ \log \frac{P(\psi_\xi^T \mid T, \theta, \xi) P(T, \theta, \xi)}{P(\psi_\xi^T, T, \xi)} - \log \frac{P(\theta, T, \xi)}{P(T, \xi)} \right] \\
&= \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} \left[ \log \frac{P(\psi_\xi^T \mid T, \theta, \xi) P(T, \theta, \xi)}{P(\psi_\xi^T \mid T, \xi) P(T, \xi)} - \log \frac{P(\theta, T, \xi)}{P(T, \xi)} \right] \\
&= \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} \left[ \log \frac{P(\psi_\xi^T \mid T, \theta, \xi)}{P(\psi_\xi^T \mid T, \xi)} \right]
\end{aligned}$$

266 Evaluating this expression requires us to compute a conditional probability with respect to  $\theta$ , which  
267 is a random variable capturing our current uncertainty over the reward network weights. We can  
268 approximate the distribution  $p(\theta)$  by sampling  $M$  weights  $\theta_1, \theta_2, \dots, \theta_M \sim p(\theta)$  and then treating  $\theta$   
269 as if it were a uniform mixture over the samples; i.e.

$$p(\theta) \approx \frac{1}{M} \sum_{i=1}^M \delta_{\theta=\theta_i},$$

270 where  $\delta_{\theta=\theta_i}$  denotes a Dirac distribution at  $\theta_i$ . Using this approximation, our mutual information  
271 becomes:

$$\begin{aligned}
I(\theta; \psi_\xi^T \mid T, \xi) &\approx \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} \left[ \log \frac{P(\psi_\xi^T \mid T, \theta, \xi)}{\frac{1}{M} \sum_{j=1}^M P(\psi_\xi^T \mid T, \theta_j, \xi)} \right] \\
&= \mathbb{E}_{\theta, \psi_\xi^T, \xi, T} \left[ \log \frac{M \cdot P(\psi_\xi^T \mid T, \theta, \xi)}{\sum_{j=1}^M P(\psi_\xi^T \mid T, \theta_j, \xi)} \right] \\
&\approx \frac{1}{M} \sum_{i=1}^M \sum_{\psi_\xi^T} P(\psi_\xi^T \mid T, \theta_i, \xi) \cdot \log \frac{M \cdot P(\psi_\xi^T \mid T, \theta_i, \xi)}{\sum_{j=1}^M P(\psi_\xi^T \mid T, \theta_j, \xi)}
\end{aligned}$$

272 **D Detailed Experiment Settings**

273 **Offline Dataset and Tasks** In real-world settings where failures are much more costly than  
274 minor failures, rollout datasets will be skewed towards higher-return rollouts. We emulate this by  
275 constraining the number of rollouts in the low reward region (return  $\leq 5$ ) to be at most 10% of the  
276 dataset.

277 **Query and Data Budgets** Preference comparison algorithms typically obtain  $n$  pairs of (query,  
278 label),  $[(\xi_{i1}, \xi_{i2}), \text{label}_i]_{i=1}^n$  for binary classification. For simpler tasks (Empty, DynObs, Lava, and  
279 Fetch), we set the query budget to be 300. That is, baselines have access to 300 (query, label)  
280 pairs,  $[(\xi_{i1}, \xi_{i2}), \text{label}_i]_{i=1}^{300}$ . IMPEC requires additional queries to precisely rank each sampled  
281 rollout within the candidate list, so we constrain it to use 150 pairs  $[(\xi_{i1}, \xi_{i2}), \text{label}_i]_{i=1}^{150}$ , and use the  
282 remaining 150 query budget to perform insertion sort for a selected subset of rollouts (decided by IG).  
283 For the harder tasks (Go to Door and Lava-Position), all algorithms are given a budget of 600 (query,  
284 label) pairs. IMPEC can access 400 data pairs, and use the remaining 200 query budget for sorting.



285 **Dataset Creation** For each task, we first train an RL agent to the expert level, saving its policies at  
 286 various timesteps. We then take 3 of its policy snapshots - an almost random policy, an expert policy,  
 287 and one in-between to generate rollouts. The number of rollouts falling within the low-reward ( $\leq 5$ )  
 288 region are controlled to take up within 10% of the dataset.

289 **IMPEC Training** We typically train an algorithm with 20 epochs. For IMPEC, we evenly divide  
 290 its query budget from epoch 1 to epoch 15, using up all queries in this period and ranking as much  
 291 uncertain rollouts as possible. We then use the remaining 5 epochs for learning the full dataset.  
 292 After sampling the initial preference pairs, IMPEC will not obtain new rollouts from the dataset, and  
 293 perform learning only by querying humans for ranking the given rollouts.

294 **Environment Observations** The environment observations are in the vector form, which contains  
 295 [agent position, agent direction, special grid info]. The special grid can be the goal/lava/door, etc.,  
 296 and its information is an encoding of its grid type, current position, and color.

297 **Ablation Studies** In the “no active learning” experiment, we turn off the active selection function,  
 298 and randomly pick rollouts from the candidate list. For “no derived prefs”, we remove the preference  
 299 derivation part of the learning. Note that the preference pairs generated during the sorting process are  
 300 still added to the dataset. Finally, “no ranking” means that the algorithm still selects preference pairs  
 301 with an information gain acquisition function, but does not maintain a preference chain (so there are  
 302 also no transitively derived preferences). This is simply the information gain algorithm of [2].

## 303 E Training Hyperparameters

304 The preference learning hyperparameters:

Hyperparameter	Value
All algorithms	
Optimizer	Adam
Learning Rate	1e-4
Weight Decay	3e-5
Batch Size	32
Temperature	0.1
Fragment Length	30
Training Epochs	20
IMPEC	
Max. Preference Chain Size	30
Stop querying at	Epoch 15
M	10

305 The PPO training hyperparameters:

Hyperparameter	Value
Training steps	500,000
Learning rate	0.0017
Gamma	0.98
Lambda	0.975
Entropy Coefficient	0.15
Batch size	64
Clip range	0.2

306 **F Ablation Studies**

307 To understand what leads to IMPEC’s performance, we experiment with removing three different  
 308 components: (1) the active learning process; (2) preference derivations; and (3) the ranking process.  
 309 The results can be found in Table 3.

310 We conduct each ablation experiment  
 311 with 25 seeds and report the number  
 312 of failed runs for each algorithm. A  
 313 run fails when its final policy’s aver-  
 314 age return  $\leq 10$ .

315 Our results suggest that there is no  
 316 one component that is responsible for  
 317 the entire gap between IMPEC and  
 318 the baseline. However, the combina-  
 319 tion of ranking and active learning can  
 320 be quite powerful: comparing the “no  
 321 derived prefs” and the baseline, the  
 322 failure rate immediately dropped by 50%.

Table 3: The failure percentages and their corresponding p-values of being significantly lower than the failure rate of the baseline

ALGORITHM-VARIANT	FAIL %	P-VALUE ( $F\% < \text{BASELINE}$ )
IMPEC	2/25	0.082
NO ACTIVE LEARNING	4/25	0.267
NO DERIVED PREFS, & NO RANKING	3/25	0.161
BASELINE	6/25	-

323 **G Graph-theoretical approach**

324 **The Preference Datasets** We visualize the preference datasets gathered by the baseline and IMPEC  
 325 on Lava-Position in Figure 5. The baseline dataset is randomly sampled at the start of the training,  
 326 while we take a snapshot of the IMPEC dataset (which is constructed iteratively) at its last training  
 327 epoch. The datasets are visualized as graphs, with each node being a unique rollout, and each edge  
 328 representing a preference label. Both IMPEC and the baseline have a query budget of 600 pairwise  
 329 queries, and IMPEC uses the first 400 of its queries to do pairwise comparisons between randomly  
 330 selected rollouts, as opposed to actively selected rollouts. The IMPEC and baseline datasets have 405  
 331 and 538 unique rollouts as well as 791 and 594 unique edges, respectively. We include several other  
 332 graph statistics in Table 4.

333 IMPEC and the baseline exhibit three notable differences in the graph properties. The first is their  
 334 clustering coefficient, which measures the degree to which nodes tend to cluster together. The  
 335 IMPEC graph has a higher clustering coefficient because of the many preferences it derives from the  
 336 ranked chain. This is relevant to the number of chains in the graph: since most nodes are connected  
 337 to IMPEC’s central cluster through some edges, it creates many more chains between two nodes  
 338 across the graph. In both graphs, we also observe that there are many chains of length 1 that are not  
 339 connected to any larger cluster. We suspect that the algorithms’ sample efficiency can be further  
 340 improved if these pairs can be meaningfully linked to the clusters.

341 Finally, we measure the graph’s efficiency, which is a metric from network science that is intended  
 342 to measure the flow of information between “communicating” nodes [11]. The assumption which  
 343 underlies the graph efficiency metric is that more distant nodes are less efficient at information  
 344 exchange. To avoid inflating the metric with the many length 1 that are not connected to the rest

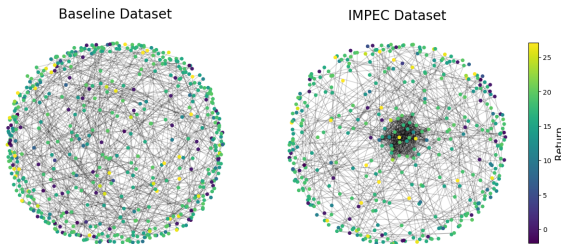


Figure 5: The preference datasets visualized as graphs. The IMPEC dataset connects more scattered rollouts through its central cluster.

	IMPEC	Baseline
Num. Nodes	405	538
Num. Edges	791	594
Cluster Coef.	0.0642	0.0013
Num. Chains	416	81
Efficiency	0.19	0.11

Table 4: Properties of the dataset graph.

345 of the graph, we compute efficiency only on the two graphs’ biggest connected components. We  
 346 empirically observe that the IMPEC dataset graph has a higher average global efficiency than the  
 347 baseline graph, which means that the average shortest path between vertex pairs in IMPEC is shorter  
 348 than for the baseline. This raises an interesting question: Is the assumption in network theory, where  
 349 the distance of nodes influences information efficiency, also applicable to preference learning? We  
 350 do not have matured results establishing connections between data connectivity in RLHF and the  
 351 training performance yet, but it would be an interesting next step for research.

## 352 H The Complete P-Value Table

Table 5: A complete list of all p-values of the algorithms performing better than the baseline for all tasks

	IMPEC	Information Gain	Volume Removal
Empty	0.10	0.82	0.37
Dynamic Obstacles	0.09	0.66	0.77
Lava	0.08	0.26	0.44
Lava Position	0.01	0.46	0.20
Fetch	0.26	0.65	0.45
Go To Door	0.37	0.55	0.39

## 353 I Baseline Comparison and Data Scaling

354 We expect that the baseline should suffer less from reward confusion if given more comparison data.  
 355 To test this hypothesis, fig. 6 shows return for the baseline with different query budgets, along with  
 356 IMPEC results from our main experiments (which used 150 preference pairs and 300 queries). We  
 357 tested the baseline with  $6\times$  the data used by IMPEC (i.e., 900 preference pairs and queries), then  
 358 gradually pushed up the amount to 1500 ( $10\times$  of IMPEC data). With  $6\times$  the data, we still see some  
 359 learning failures, and the seeds’ standard deviation decreases to IMPEC’s level only when we use  
 360  $10\times$  data.

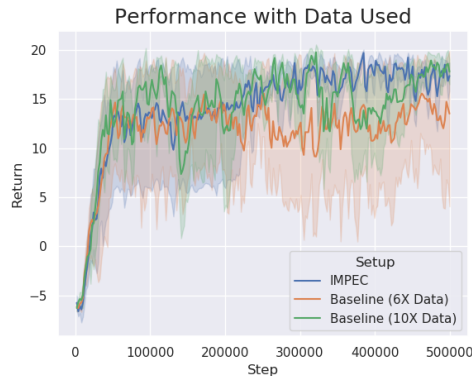


Figure 6: The learning curves for IMPEC and baseline with different amounts of data.

361 **J Learning Curves for All Tasks**

