

---

# Bayesian Hypothesis Testing Policy Regularization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        In reinforcement learning (RL), sparse feedback makes it difficult to target long-  
2        term outcomes, often resulting in high-variance policies. Real-world interventions  
3        instead rely on prior study data, expert input, or short-term proxies to guide explo-  
4        ration. In this work, we propose *Bayesian Hypothesis Testing Policy Regularization*  
5        (BHTPR), a method that integrates a previously-learned policy with a policy learned  
6        online. BHTPR uses Bayesian hypothesis testing to determine, state by state, when  
7        to transfer the prior policy and when to rely on online learning.

## 8 1 Introduction

9 In RL settings with sparse or delayed rewards, exploration can be costly or risky—particularly in  
10 mobile health studies, where outcomes can only be observed infrequently or at the end of the study.  
11 Excessive exploration can overburden participants, causing disengagement. To reduce variance and  
12 speed up learning, policies guided by imperfect mediators, proximal outcomes, or previous study data  
13 can help guide the agent’s exploration. For example, in the mobile health algorithm “HeartSteps”,  
14 the RL agent targeted a “short-term, measurable behavioral or psychosocial effect through which  
15 that component is hypothesized to mediate desired distal health outcomes” [3]. We can view these  
16 proximal outcomes in terms of their induced policy, guiding the agent towards the distal outcome.

17 Policy regularization is commonly used in offline RL to reduce the learned policy’s deviation from  
18 the behavior policy [13]. In this setting, policy regularization reduces bias in the learned policy  
19 caused by optimism when extrapolating beyond the observed data [4]. In this paper, we instead apply  
20 policy regularization to online learning with sparse, distal rewards. In this setting, rather than prevent  
21 optimism, policy regularization guides learning using previous study data.

22 In this paper, we introduce Bayesian Hypothesis Testing Policy Regularization (BHTPR), an online  
23 method using Bayesian hypothesis testing to leverage previous study data. The method allows us to  
24 utilize this data even when the transition dynamics in the previous study environment differ from  
25 those in the environment of interest. In BHTPR, the agent selects the action at each timestep from  
26 either the policy calculated from the previous data or from the policy learned online. The probability  
27 of acting according to each policy is determined by a Bayesian hypothesis test. Because we do not  
28 directly combine Q-functions, our method allows combining data even when the outcomes from the  
29 two environments are on different scales. Furthermore, by falling back on the prior policy in states  
30 where the transition dynamics match, the agent avoids unnecessary exploration.

31 We demonstrate the performance of BHTPR across a range of experiments varying episode length,  
32 quality of prior data, environment stochasticity, and other factors. BHTPR demonstrates good  
33 performance both on a simple tabular example and on a more realistic mobile health setting. We  
34 compare BHTPR to (1) standard epsilon-greedy exploration, and (2) using the previous study data  
35 as a prior on the transition function. We find that our method performs similarly to or better than  
36 baselines in variations in the setting and environment and is robust to parameter choices.

## 37 2 Related Works

38 **Approaches for Sparse, Distal Rewards** The problem of pursuing a distal outcome in the face of  
39 sparse feedback is a core problem of RL that has been addressed by many areas of research. Reward  
40 design focuses on providing structured signals that guide behavior and accelerate learning [9, 10, 2].  
41 Reward design is a key challenge for mobile health studies [12]. Reward shaping augments the  
42 environment’s reward function with additional feedback to guide the agent to learn more quickly  
43 without altering the optimal policy [6]. Safe RL focuses on problems “in which it is important to  
44 ensure reasonable system performance and/or respect safety constraints during the learning and/or  
45 deployment process” [1]. Finally, transfer learning leverages experience gained from similar tasks  
46 to improve the learning of a novel task [11, 5]. In this work, we modify exploration using external  
47 knowledge. We build on prior works by using Bayesian hypothesis testing to reflect the inductive  
48 bias that the external knowledge is correct at some states and incorrect at others.

49 **Policy Regularization** Policy regularization is commonly used to avoid a policy deviating too  
50 far from its prior version or from a reference policy, improving stability, sample efficiency, and  
51 generalization. In offline settings, policy regularization ensures the learned policy does not deviate  
52 far from the behavior policy, reducing extrapolation error [13]. In online settings, methods such as  
53 Trust Region Policy Optimization (TRPO) constrain successive policy updates using KL divergence  
54 penalties[8]. Our method can be viewed as policy regularization, as it constrains the learned policy  
55 to use a policy derived from prior study data in certain states. Furthermore, it can be framed within  
56 BRAC, as discussed in Sec. 6.

57 **Bayesian Hypothesis Testing** We are aware of one other work that explicitly incorporates Bayesian  
58 hypothesis testing (BHT) into an RL algorithm. Zhang et al. [14] uses BHT to interpolate between  
59 MDP and bandit algorithms when the true nature of the environment is not known. In contrast, our

60 method assumes that the true environment is an MDP and the role of BHT is to learn the states in  
61 which the given policy is correct.

## 62 **3 Methods**

### 63 **3.1 Setting and Notation**

64 We consider an episodic Markov decision process (MDP) setting with, with states  $s \in \mathcal{S}$ , actions  
65  $a \in \mathcal{A}$ , reward function  $R(s, a)$  and transition function  $T(s, a, \cdot)$ . In the experiments below, we  
66 assume  $R(s, a)$  is known. We define  $\mathcal{D}$  to be the transition data collected online. We use this data  
67 and the reward function to learn a policy  $\pi$  and an MLE transition matrix  $\hat{T}$ . We define  $\mathcal{D}_{prior}$   
68 as the transition data collected from some prior study. We compute the optimal policy based on  
69 this transition data and the know reward function  $R(s, a)$ , and call this policy  $\pi_{prior}$ . We compute  
70  $\hat{T}_{prior}$  as the MLE transition matrix using  $\mathcal{D}_{prior}$ . The transition function generated from  $\mathcal{D}_{prior}$   
71 may match that of the new environment or it may differ in certain states. Consequently, we aim to  
72 rapidly differentiate between states where  $\pi_{prior}$  can be utilized and states where a new policy must  
73 be learned online.

### 74 **3.2 Algorithm Definition**

75 BHTPR learns weights  $w(s)$  that reflects the agent’s belief about whether the transition dynamics in  
76 the current environment match those from the prior study in each state  $s$ . Then  $w(s)$  is used within  
77 the algorithm to decide whether to select an action in state  $s$  according to  $\pi$  (the policy being learned  
78 online) or  $\pi_{prior}$  (the policy computed from the prior study data).

#### 79 **3.2.1 Update rule for $w(s)$ using Bayesian hypothesis testing**

80 BHTPR assumes access to data  $\mathcal{D}_{prior}$  that comes from an environment with the same state space, but  
81 potentially different transition dynamics in some states. In order to differentiate between the states in  
82 which the transition dynamics are the same and those that differ, BHTPR learns a state-specific weight  
83  $w(s) \in [0, 1]$  that represents the probability that the transition dynamics in the current environment  
84 match those implied by the prior data. This is formalized via Bayesian hypothesis testing, which  
85 updates  $w(s)$  as more transition data are collected.

86 In each state, we test the null hypothesis that transitions match the prior environment against the  
87 alternative hypothesis that they do not. Using observed transition data, we update  $w(s)$  to reflect the  
88 posterior probability that the prior model is correct in that state. We use Bayes’ rule to compute this  
89 posterior probability, based on the likelihood of transition data under the prior and learned models.  
90 The update equation and further details are provided in Appx. A.

#### 91 **3.2.2 BHTPR Algorithm: Using $w(s)$ to choose between policies**

92 At every step of episode  $e$ , the online algorithm chooses the policy from which to select its action  
93 based on a draw from a Bernoulli trial, where the probability of acting according to  $\pi_{prior}(s)$  is  
94  $w_e(s)$ . If the resulting value  $b \sim \text{Bernoulli}(w_e(s))$  equals 1, the agent selects its action based on  
95  $\pi_{prior}(s)$ , otherwise it selects an action from the current learned policy  $\pi(s)$ . Hence the agent takes  
96 actions based on the previous study data in proportion to the probability that the observed data was  
97 generated from the previous study transition function. This procedure is detailed in Algorithm 1.

---

**Algorithm 1** Bayesian hypothesis testing policy regularization

---

```
1: Input: Previous study data  $\mathcal{D}_{prior}$ 
2: Calculate  $\hat{T}_{prior}, \pi_{prior}$  using  $\mathcal{D}_{prior}$ ; Initialize  $\hat{T}(s, a, \cdot)$ , weight  $w_0(s)$ , policy  $\pi$ 
3: for each episode  $e$  do:
4:   Initialize episode data  $\mathcal{D}_e(s) = \{\} \forall s$ 
5:   Initialize state  $s$  from starting state distribution
6:   while  $s$  not terminal do
7:     Draw  $b \sim \text{Bernoulli}(w_e(s))$ 
8:     if  $b = 1$  then
9:       Choose action  $a \sim \pi_{prior}(s)$ 
10:    else if  $b = 0$  then
11:      Choose action  $a$  epsilon-greedily from  $\pi(s)$ 
12:    end if
13:    Take action  $a$ . Observe next state  $s'$ . Update data  $\mathcal{D}_e(s) = \mathcal{D}_e(s) \cup (s, a, s')$ .
14:    Update  $\hat{T}, \pi$ 
15:     $s \leftarrow s'$ 
16:  end while
17:  Update  $w$  using Bayesian hypothesis testing (Eq. 1)
18: end for
```

---

## 98 4 Experiments

99 We demonstrate the ability of BHTPR to leverage previous study data to speed up learning in settings  
100 with sparse, distal rewards. We show how the algorithm can outperform a standard epsilon-greedy  
101 approach as well as an epsilon-greedy approach using the previous study data directly as a Dirichlet  
102 prior on the transition function. Additional experimental details are in Appx. C.

### 103 4.1 Environments

104 We illustrate the performance of BHTPR on experiments using two environments. The first is a grid  
105 in which an agent must navigate from one corner to the opposite corner. This simple example allows  
106 us to manipulate different aspects of the environment and observe the effect on the performance of  
107 BHTPR. The second environment is the “Chainworld” from Nofshin et al. [7]. This reflects a realistic  
108 depiction of the challenges faced in a mobile health setting.

109 **Gridworld** The state space of this environment is a 7x7 grid with four actions (up/down/left/right).  
110 The agent starts at the lower-left of the grid and the episode concludes when either the agent reaches  
111 the upper-right corner or after a fixed number of steps. The agent receives a reward of -1 per step  
112 until reaching the goal. When the agent reaches the goal, it receives a reward of 100 and the episode  
113 terminates. This is depicted in Fig. 4a in Appx. B.

114 **Chainworld** A more realistic mobile health example, the “Chainworld” from Nofshin et al. [7],  
115 applies behavioral science literature to model human-AI interaction when an AI assists a human in  
116 a frictionful task such as adhering to a physical therapy program. Both the human and the AI are  
117 represented by the MDPs, where the AI MDP’s actions affect the human’s MDP. The human MDP  
118 represents progress towards the goal and the AI’s MDP uses the human’s MDP state and action as its  
119 state space. Further details are in Appx. B.

## 120 5 Results

121 As we discuss below, BHTPR has good performance across variations in the gridworld and chainworld  
122 environments. In the chainworld, the results are greatly shaped by the presence of the absorbing state.

### 123 5.1 Impact of an Absorbing State on Exploration

124 The main element distinguishing the chainworld environment is the consequence of exploration. In  
125 contrast to the gridworld, where unnecessary exploration increases the number of steps to reach the

126 goal, exploration in the chainworld can result in disengagement after which the agent can never reach  
 127 the goal state. Hence exploration is incredibly costly and the best algorithm in this environment is  
 128 one that minimizes unnecessary exploration. BHTPR can reduce exploration because, once the agent  
 129 has learned that  $\pi_{prox}$  is correct, it can stop exploring in that state. In terms of Algorithm 1, this  
 130 means acting epsilon-greedily only if  $b = 0$  (acting according to the policy  $\pi$  learned online, line 11)  
 131 and not exploring when acting according to  $\pi_{prior}$  (line 9).

132 Fig. 5 in Appx. D highlights the impact of exploration on the performance of BHTPR and the  
 133 baselines in the chainworld. In the right column, equal exploration is enforced for all methods. Here  
 134 BHTPR performs similarly to the prior on  $T$ . This holds in both the case when there is a large  
 135 negative reward for disengagement (bottom row) and when there is not (top row). Epsilon-greedy is  
 136 slower to learn without the large disengagement reward but performs similarly to the other two in the  
 137 presence of a large disengagement reward. In contrast, when BHTPR is not forced to explore equally  
 138 with the baselines, it outperforms the baselines, as seen in the left column of Fig. 5.

139 For the sake of comparison between methods, in the gridworld example, we fix the amount of  
 140 exploration to be equal across BHTPR and the baselines. In the chainworld examples, however, since  
 141 performance is driven by exploration, BHTPR is implemented in the results that follow with *unequal*  
 142 *exploration* in the chainworld environment (i.e. No exploration when  $b = 1$ ).

143 **5.2 BHTPR can learn over short episodes.**

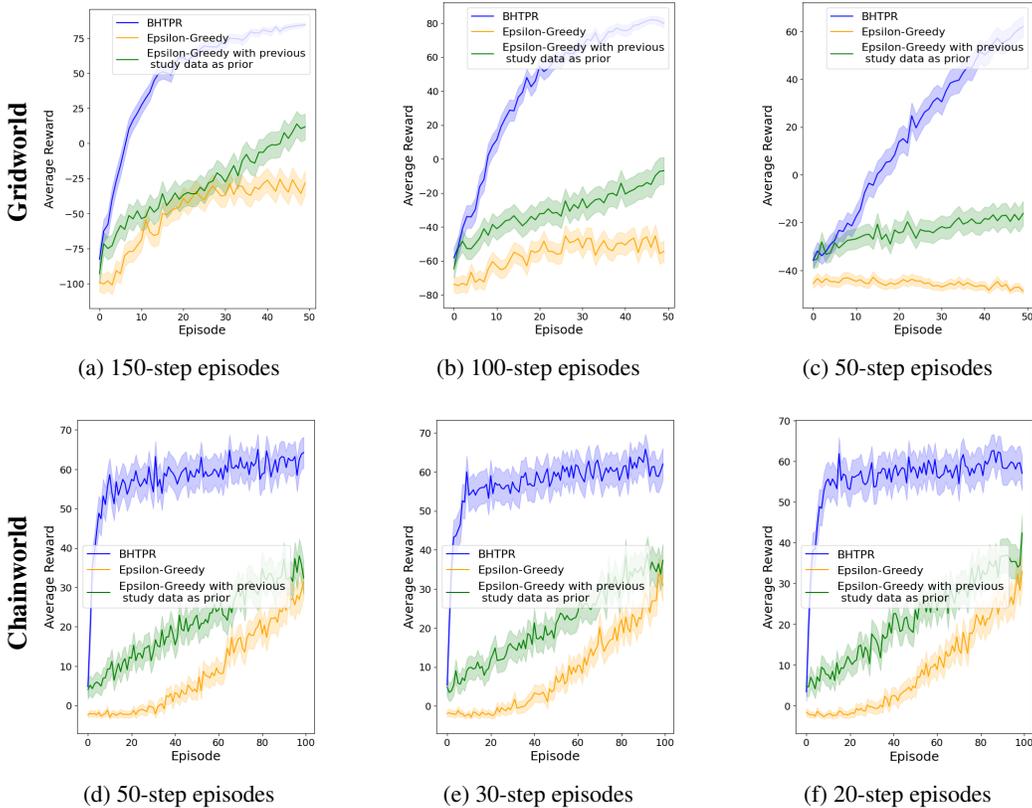


Figure 1: **Episode length.** BHTPR learns in gridworld even with short episodes, unlike the baselines. In chainworld, episode length has limited effect due to disengagement.

144 **Gridworld** By incorporating correct actions learned from the previous study data, BHTPR is able  
 145 to learn over short episodes. This is the case even when episodes are too short for the baselines to  
 146 learn. Fig. 1 illustrates the performance of our method compared to the baselines in environments  
 147 of different episode lengths. In the gridworld, all three methods are able to learn when the episodes  
 148 are sufficiently long (Fig. 1a), however, when the episodes are short (Fig. 1c), the epsilon-greedy  
 149 agent cannot learn at all and epsilon-greedy agent with the prior study data as a prior on  $T$  learns

150 very slowly. BHTPR is able to learn with a greatly reduced number of steps per episode compared to  
 151 standard epsilon-greedy learning because it quickly learns the states in which  $\pi_{prior}$  is the correct  
 152 policy and does not have to spend as many steps in the episode taking incorrect actions.<sup>1</sup> To confirm  
 153 this intuition, see Fig. 6 in Appx. D. Here, we mix actions from the true optimal policy with the  
 154 learned policy for a range of fixed percentages (25, 50, and 75 percent). The figure shows that as the  
 155 proportion of correct actions increases, the performance approaches that of BHTPR, demonstrating  
 156 that the reason our method can learn with fewer episode steps is due to incorporating correct actions.

157 **Chainworld** In contrast to the gridworld, the chainworld’s absorbing state limits the impact of the  
 158 number of episode steps on the performance of each method. This is illustrated in Fig. 1(d-f), where  
 159 we see little impact of the change in maximum steps per episode across the plots in the row. Once the  
 160 disengagement state is reached, no further learning is possible, even if the agent has not reached the  
 161 maximum number of steps in the episode.

### 162 5.3 BHTPR performs well regardless of the correctness of the prior data.

163 When prior study data are used to form the policy for a new study, this policy will often be incorrect  
 164 in certain regions of the state space, for example because the new study is performed on a population  
 165 that differs from that of the previous study. Since we do not know the number or identity of these  
 166 states, BHTPR must perform well regardless of the number of states for which the prior study data  
 167 are incorrect for the new setting. In Fig. 2, we vary the size of the region in which the previous study  
 168 data matches the current environment. For the states in which the data are corrupted, the environment  
 169 for the prior study data is randomly generated by selecting the next state for the transition function  
 170 uniformly at random. In the first column (Figs. 2a and 2d),  $\mathcal{D}_{prior}$  matches the environment perfectly,  
 171 whereas in the next two columns, the prior study data’s environment is increasingly different from the  
 172 target environment. In the gridworld examples, when the prior study data is perfect, using it directly  
 173 as a prior on the transition function performs best, however BHTPR is very close in performance,  
 174 greatly outperforming the epsilon-greedy strategy. Note that this is not reflected in the chainworld  
 175 example (Fig. 2d) because we do not fix exploration to be equal across methods. Hence, BHTPR  
 176 performs better even when the data is perfect because of the lower amount of exploration, otherwise  
 177 the baseline in which the previous study data serves as a prior would be favored in this case as well.  
 178 This highlights how much exploration is driving the performance in the chainworld. When the data  
 179 are partially correct, our method outperforms the two baselines. Note, however, that with greater  
 180 regions of incorrect data, the exploration parameter needs to be decayed more slowly otherwise  
 181 BHTPR will fail to learn the optimal policy.

182 The performance of BHTPR across the range of data corruption schemes shown here is related to our  
 183 inductive bias that the data are either entirely correct or entirely wrong in each state. In contrast, in a  
 184 setting where there is a distribution shift between the prior study environment and the environment of  
 185 interest, we would expect the baseline in which the previous study data forms a prior to perform well.  
 186 We confirm this intuition in Fig. 12. In these experiments, we generate an incorrect transition function  
 187 in the corrupted states by drawing a next state uniformly at random, as before, but then instead of  
 188 using this transition function directly, we mix the incorrect and correct transition functions together  
 189 in a range of proportions to create a distribution shift between the previous study environment and  
 190 the new environment. When the prior is close to the target environment (Fig. 12a), the Bayesian prior  
 191 baseline performs slightly better; as the mismatch increases, BHTPR outperforms.

### 192 5.4 BHTPR performs well in environments of varying levels of stochasticity.

193 Varying the stochasticity of the gridworld environment, the relative performance of BHTPR is greater  
 194 in a more deterministic environment, as demonstrated in Fig. 3. This behavior reflects the setup of  
 195 experiments where the amount of prior study data is fixed across examples, so the maximum prior  
 196 magnitude  $T_{prior}(s, a, \cdot)$  is highest for more deterministic environments. To confirm this, we fix the  
 197 distribution of prior  $T_{prior}(s, a, \cdot)$ , but vary the magnitude in Fig. 7 (Appx. D). Given the correct  
 198 prior magnitude, the baseline Bayesian method can outperform; however, our method avoids this  
 199 parameter tuning and, as demonstrated in the next section, is robust across initializations of  $w(s)$ .

<sup>1</sup>Exploration is equalized across methods in the Gridworld example, however fewer wrong actions are still taken by BHTPR because in the states  $\pi_{prior}$  is correct, the agent rapidly learns to use this. The baselines take more incorrect actions while learning the policy online.

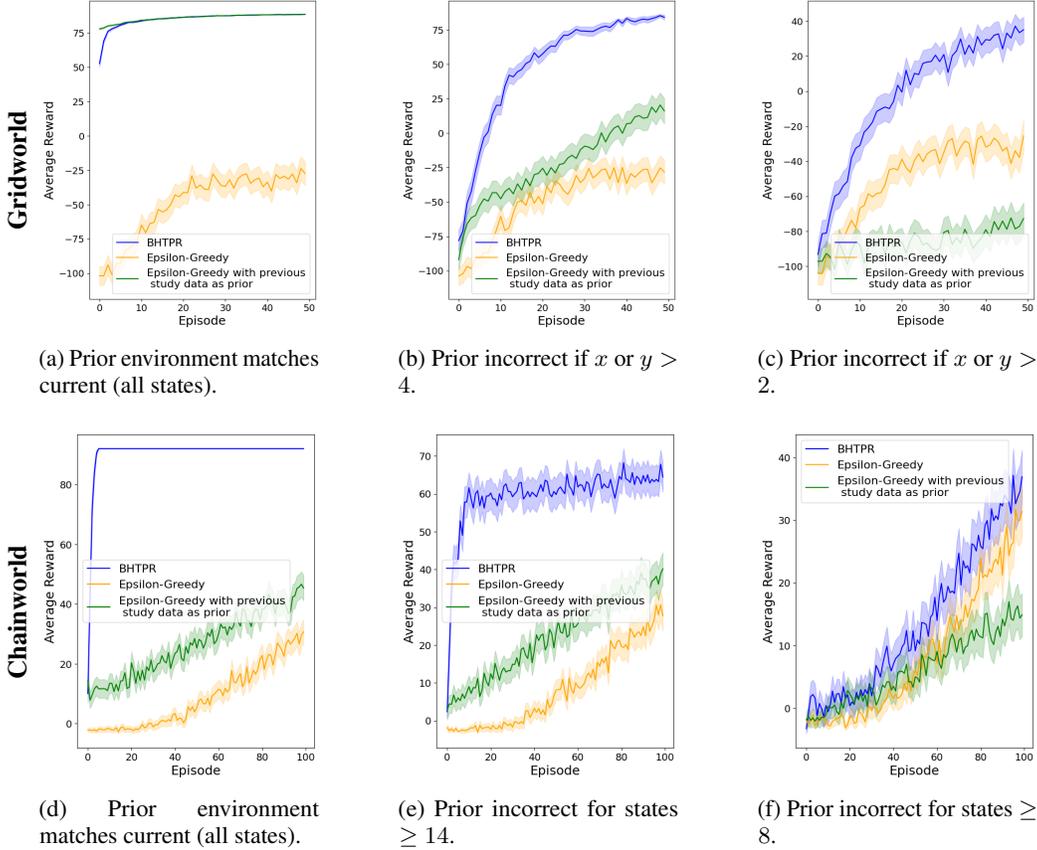


Figure 2: **Correctness of prior study data.** BHTPR performs best when prior and current environments match (a,d), but still outperforms or matches baselines when prior data are partially incorrect.

200 In the chainworld environment, BHTPR performs well across variations in the human’s probability of  
 201 moving, but all methods perform best in the deterministic environment. In the leftmost panel of Fig. 3,  
 202 the human moves with 100% probability if they take an action, and this probability decreases from  
 203 left to right. This results in the probability of the human disengaging if they take action 0 decreasing  
 204 from left to right (d-f). Moving from plot (d) to (e) to (f), the agent will encounter the disengagement  
 205 state less frequently and consequently receive less feedback about the possibility of disengagement  
 206 when taking action 0. This makes it harder for the agent to learn from random exploration, and  
 207 consequently the benefit of injecting outside information is particularly strong in this case.

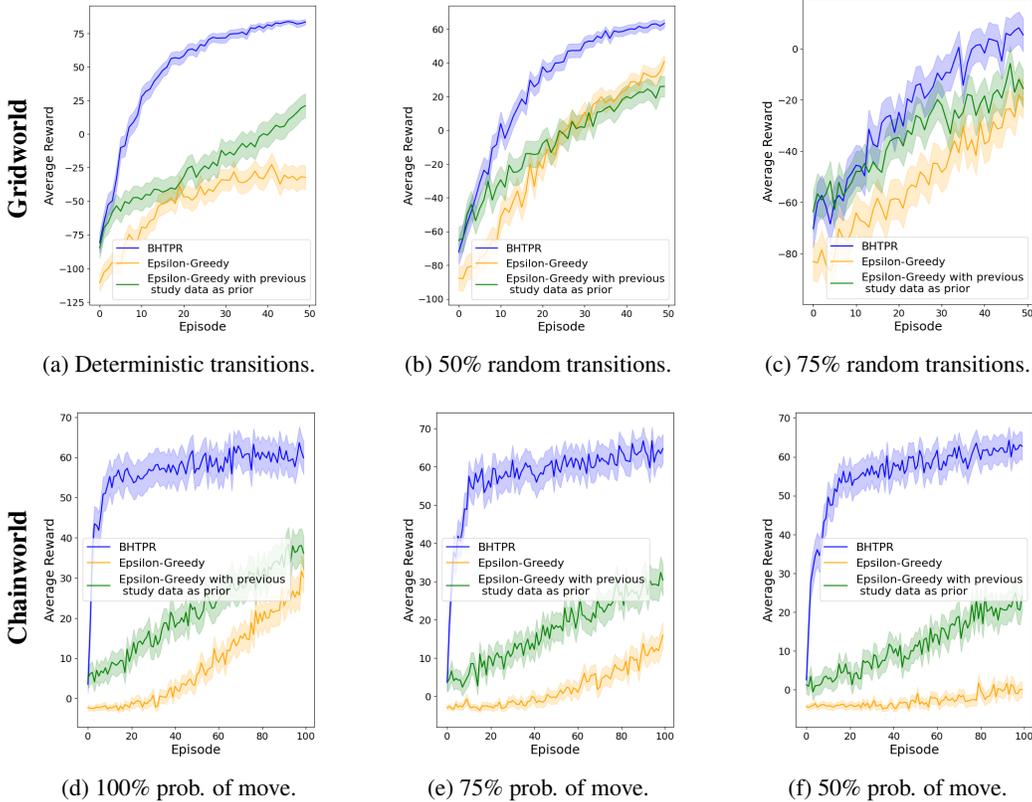


Figure 3: **Stochasticity.** In gridworld, BHTPR’s advantage is largest in deterministic environments. In chainworld, all methods perform better when the human behaves deterministically.

208 **6 Discussion**

209 **Connection to other policy regularization methods** As discussed in Sec. 2, BHTPR can be viewed  
 210 as a form of policy regularization. Wu et al. [13] define a framework, BRAC, to describe different  
 211 methods of policy regularization for off-policy RL, and we can view BHTPR as an instance of this  
 212 framework. BRAC describes policy regularization methods in terms of a penalty on the Q-value  
 213 objective and/or the policy objective based on how the learned policy diverges from the behavior  
 214 policy, encouraging similarity between the learned and behavior policies. Similarly, our method  
 215 also encourages the learned policy to be more similar to  $\pi_{prior}$ , the policy learned from the data of  
 216 the previous study. Bayesian hypothesis testing provides a principled way to set the strength of the  
 217 penalty— here the probability of acting according to  $\pi_{prox}$ .

218 **Limitations** The main limitation of this method is that it is currently framed in a tabular setting.  
 219 While it would be easy to use a weight  $w(s)$  to interpolate between learned and prior policies where  
 220 both are defined in a continuous state/action setting, extending the calculation for  $w$  itself to a function  
 221 approximation setting requires further work. We also note that compared to the baselines methods,  
 222 BHTPR introduces additional computational overhead in terms of updating weight  $w(s)$  for all states  
 223 after each episode. Finally, BHTPR does not account for the transfer of policies between settings  
 224 where the state spaces are not identical.

225 **7 Conclusion**

226 Learning in a setting with sparse, distal rewards— a situation commonly encountered in mobile  
 227 health— provides unique challenges. Policies crafted from previous study data or by experts can guide  
 228 exploration, leading the agent to learn more quickly. Our method uses Bayesian hypothesis testing to  
 229 incorporate prior study data in a principled way and speeds up learning across a range of settings.

## 230 References

- 231 [1] Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning.  
232 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 233 [2] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J Russell, and Anca Dragan. Inverse  
234 reward design. *Advances in neural information processing systems*, 30, 2017.
- 235 [3] Predrag Klasnja, Shawna Smith, Nicholas J Seewald, Andy Lee, Kelly Hall, Brook Luers,  
236 Eric B Hekler, and Susan A Murphy. Efficacy of contextually tailored suggestions for physical  
237 activity: a micro-randomized optimization trial of heartsteps. *Annals of Behavioral Medicine*,  
238 53(6):573–582, 2019.
- 239 [4] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning  
240 for offline reinforcement learning. *Advances in neural information processing systems*, 33:  
241 1179–1191, 2020.
- 242 [5] Alessandro Lazaric. Transfer in reinforcement learning: a framework and a survey. In *Rein-*  
243 *forcement Learning: State-of-the-Art*, pages 143–173. Springer, 2012.
- 244 [6] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transforma-  
245 tions: Theory and application to reward shaping. In *Icml*, volume 99, pages 278–287. Citeseer,  
246 1999.
- 247 [7] Eura Nofshin, Siddharth Swaroop, Weiwei Pan, Susan Murphy, and Finale Doshi-Velez. Rein-  
248 forcement learning interventions on boundedly rational human agents in frictionful tasks. In  
249 *Proceedings of the... International Joint Conference on Autonomous Agents and Multiagent Sys-*  
250 *tems: AAMAS. International Joint Conference on Autonomous Agents and Multiagent Systems*,  
251 volume 2024, page 1482, 2024.
- 252 [8] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust  
253 region policy optimization. In *International conference on machine learning*, pages 1889–1897.  
254 PMLR, 2015.
- 255 [9] Satinder Singh, Richard L Lewis, and Andrew G Barto. Where do rewards come from. In  
256 *Proceedings of the annual conference of the cognitive science society*, pages 2601–2606.  
257 Cognitive Science Society, 2009.
- 258 [10] Henry Sowerby, Zhiyuan Zhou, and Michael L Littman. Designing rewards for fast learning.  
259 *arXiv preprint arXiv:2205.15400*, 2022.
- 260 [11] Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A  
261 survey. *Journal of Machine Learning Research*, 10(7), 2009.
- 262 [12] Anna L Trella, Kelly W Zhang, Inbal Nahum-Shani, Vivek Shetty, Finale Doshi-Velez, and  
263 Susan A Murphy. Reward design for an online reinforcement learning algorithm supporting  
264 oral self-care. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,  
265 pages 15724–15730, 2023.
- 266 [13] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
267 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 268 [14] Kelly W Zhang, Omer Gottesman, and Finale Doshi-Velez. A bayesian approach to learning  
269 bandit structure in markov decision processes. *arXiv preprint arXiv:2208.00250*, 2022.

## 270 A Bayesian Hypothesis Testing Details

271  $w(s)$  represents the posterior probability that the observed transition data was generated by  $\hat{T}_{prior}$ ,  
272 the current estimate of  $T_{prior}$  as opposed to  $\hat{T}$ . We formalize this idea using Bayesian hypothesis  
273 testing, with null an alternative hypotheses defined as follows.

274 **Null hypothesis,  $H_0$**  The transition data collected online was generated by the transition function  
 275 implied by the prior study data, i.e.  $T(s, a, \cdot) = T_{prior}(s, a, \cdot)$  for all actions  $a \in \mathcal{A}$  in a given state  
 276  $s$ .

277 **Alternative hypothesis,  $H_1$**  The transition data collected online was not generated by the transition  
 278 function implied by the prior study data; therefore, we calculate the likelihood of the transition data  
 279 under the transition function being learned online.

280 **Per-episode update for  $w(s)$**  Using the null and alternative hypotheses above, we set  $w(s)$  equal to  
 281 the posterior probability of the null hypothesis, i.e. the probability that the prior study data is correct  
 282 in state  $s$ .

$$P(H_0|\mathcal{D}) = \frac{P(\mathcal{D}|H_0)P(H_0)}{P(\mathcal{D})} = \frac{P(\mathcal{D}|H_0)P(H_0)}{P(\mathcal{D}|H_0)P(H_0) + P(\mathcal{D}|H_1)P(H_1)}$$

283 Letting the initial weight  $w_0(s)$  be the initial  $P(H_0)$ , we calculate the likelihood of the data using the  
 284 transition function calculated from the prior data,  $\hat{T}_{prior}$ , and the transition function learned online,  
 285  $\hat{T}$ . Then after every episode  $e$  the weight  $w_e(s)$  is updated according to Eq. 1. Let  $\mathcal{D}_e(s)$  be the set  
 286 of all transition data  $\{(s, a, s')\}$  collected in episode  $e$  starting from state  $s$ .

$$w_{e+1}(s) = \frac{\underbrace{\left[ \prod_{(s,a,s') \in \mathcal{D}_e(s)} T_{prior}(s, a, s') \right]}_{\text{likelihood of data under } H_0} \underbrace{w_e(s)}_{P(H_0)}}{\underbrace{\left[ \prod_{(s,a,s') \in \mathcal{D}_e(s)} T_{prior}(s, a, s') \right]}_{\text{likelihood of data under } H_0} \underbrace{w_e(s)}_{P(H_0)} + \underbrace{\left[ \prod_{(s,a,s') \in \mathcal{D}_e(s)} T(s, a, s') \right]}_{\text{likelihood of data under } H_1} \underbrace{(1 - w_e(s))}_{P(H_1)}} \quad (1)$$

287 If we further assume that there are regions of the state space in which the current environment matches  
 288 the prior study environment and others in which it does not, we can improve upon this method by  
 289 applying smoothing to values of  $w_e(s)$  across steps in each episode so that states that are close  
 290 together have similar weights. In the examples that follow, we apply an exponential moving average  
 291 over weights at each timestep of an episode.

## 292 B Environment Details

### 293 B.1 Environment Visualizations

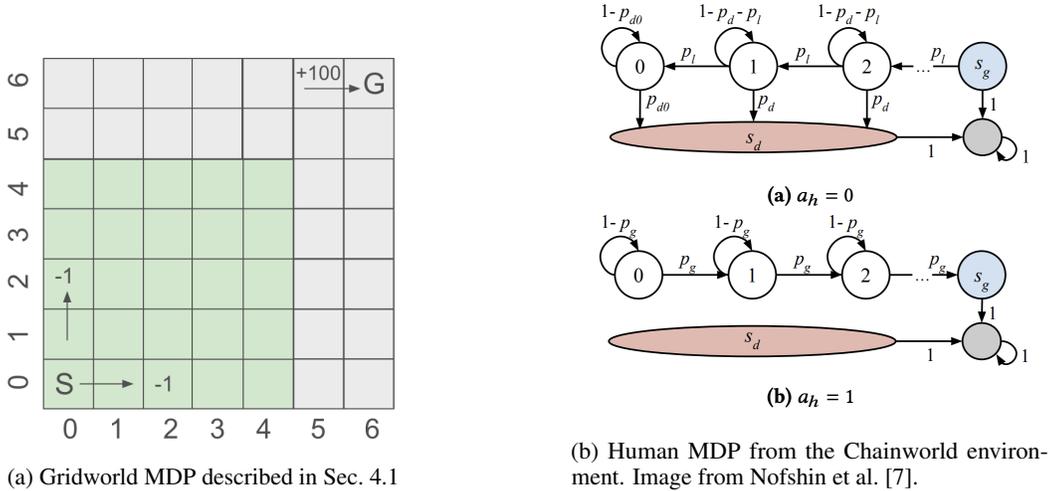


Figure 4: Two environments used in Sec. 4

## 294 B.2 Chainworld Additional Details

295 The human is modeled by MDP  $M_h = \langle S_h, A_h, T_h, R_h, \gamma_h \rangle$ , representing progress towards the  
296 task goal. Each state reflects the amount of progress the human has made towards that goal. If the  
297 human takes action 0 (Fig. 4b(a)), they may stay in the same state, lose progress towards the goal  
298 (move back one state), or disengage (enter state  $s_d$ ). If instead they take action 1 (Fig. 4b(b)), they  
299 will either remain in the same state or progress one state closer to the goal. The human takes the  
300 optimal action according to their MDP. If the human disengages, they cannot re-engage.

301 The AI is represented by MDP  $M_{AI} = \langle S_{AI}, A_{AI}, T_{AI}, R_{AI}, \gamma_{AI} \rangle$ .<sup>2</sup> The AI assists the human,  
302 who is a myopic decision-maker (low  $\gamma_h$ ), at reaching the goal state. To achieve this goal, the AI  
303 has two possible actions. It can send an intervention that temporarily increases  $\gamma_h$  (action 1) or do  
304 nothing (action 0). Its two-dimensional state space is represented by  $S_{AI} = [s_h, a_h]$ , the human’s  
305 state and previous action.

306 In the experiments below, we consider the AI agent’s MDP. Parameters for the human and AI MDPs  
307 were chosen such that the agent can reach the goal (i.e. the right combination of agent actions can  
308 lead the human to the goal), and such that  $T_{AI}$  differs for the AI’s two actions (e.g. if the human is  
309 not sufficiently myopic, he will take action 1 regardless of AI interaction and BHTPR will have no  
310 impact). The environment has 10 human states. Please see Nofshin et al. [7] for more details on the  
311 MDP construction.

312 One natural way to construct a reward function in the presence of an undesirable absorbing state is  
313 with a large negative reward in this state. In this case, the agent receives the large negative reward  
314 upon encountering the absorbing state and correctly learns to avoid taking that action again in the  
315 given state. Consequently, simple epsilon-greedy methods can learn in this setting. Since immediate  
316 feedback is available in the form of a large negative reward every time the agent encounters the  
317 disengagement state, this implementation of the chainworld does not represent a sparse, distal reward  
318 unless the action probabilities are set such that the disengagement state is not frequently encountered.  
319 To reflect a sparse rewards environment that we are primarily concerned with, results below do not  
320 use a reward function with a large negative reward in the absorbing state (instead  $R(s, a) = -1$  in  
321 the absorbing state, as in all other states), however results for this variation are available in Appx. E.

## 322 C Experiment Details

323 In the experiments, we vary the following elements: stochasticity of the transitions, correctness of  
324 the prior study data, and the maximum number of steps per episode. By default in the gridworld,  
325 unless otherwise stated, the previous study data are correct for states in which the x and y coordinates  
326 are less than 5 (illustrated in green in Fig. 4a), and the maximum episode length is 150 steps. In the  
327 chainworld examples, previous study data are correct for states numbered less than or equal to 13 and  
328 the maximum episode length is 50. In both, the weights  $w_0(s)$  in the algorithm are initialized to 0.5  
329 for all states and transitions are deterministic by default.

---

<sup>2</sup>Our algorithm considers an episodic setting where  $\gamma_{AI} = 1$ .

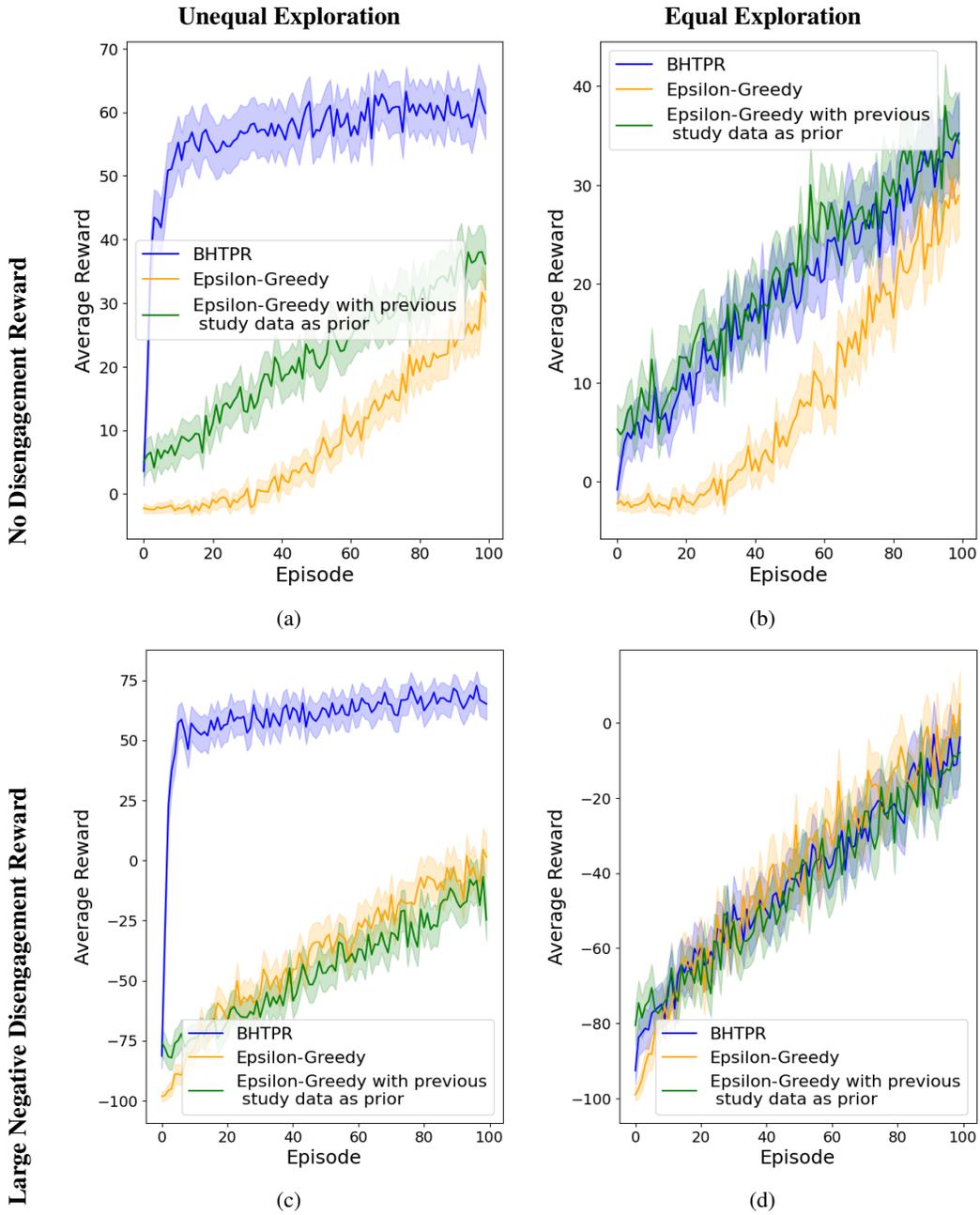


Figure 5: Effect of varying disengagement reward and exploration in the chainworld environment.

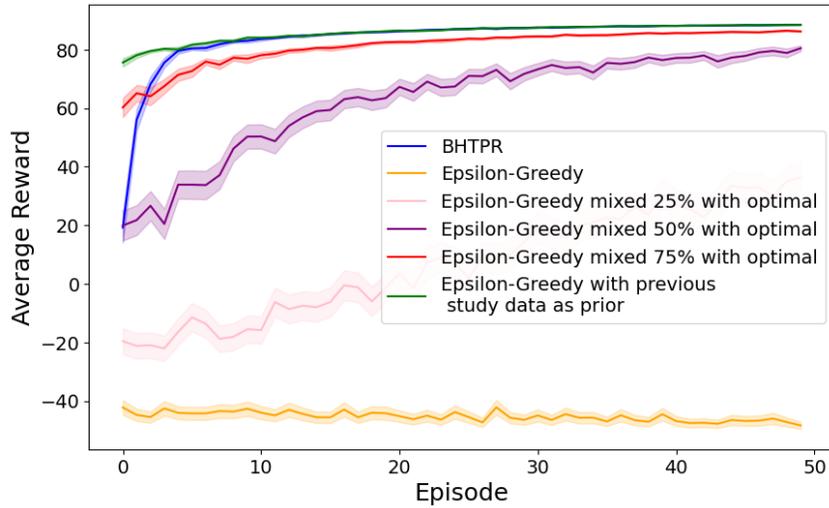


Figure 6: As the percentage of optimal actions mixed with an epsilon-greedy increases, performance approaches BHTPR. Plot from Grid MDP environment, previous study data matches current environment in all states.

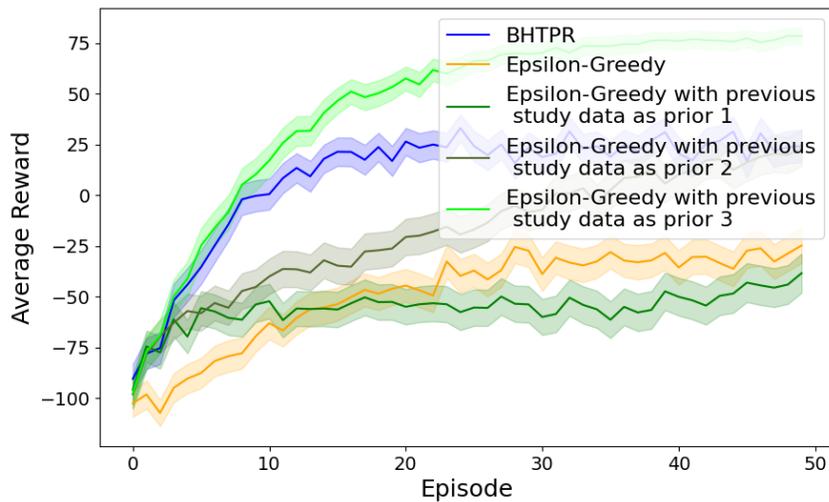


Figure 7: Using the previous study data as a prior on  $T$  can outperform BHTPR for the correct magnitude of prior, however BHTPR avoids tuning the magnitude of the prior. Grid MDP environment.

331 BHTPR is robust to a wide range of initial values for  $w_0(s)$  (Fig. 8), though performance can improve  
 332 slightly with expert-informed initializations.

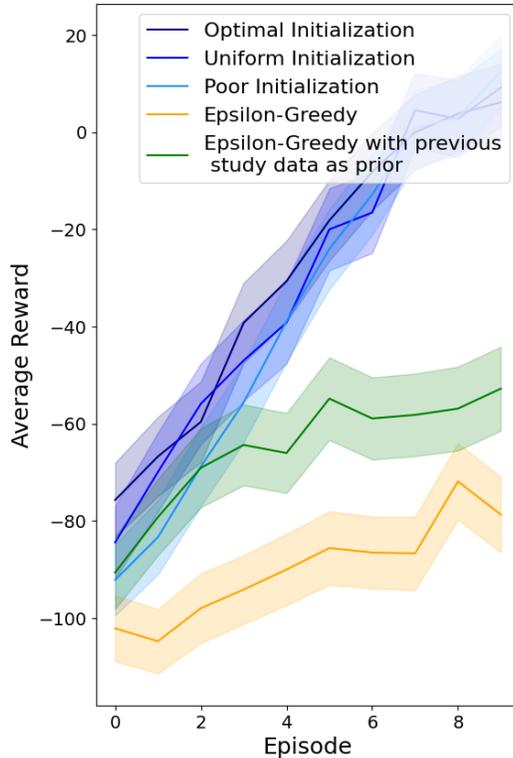
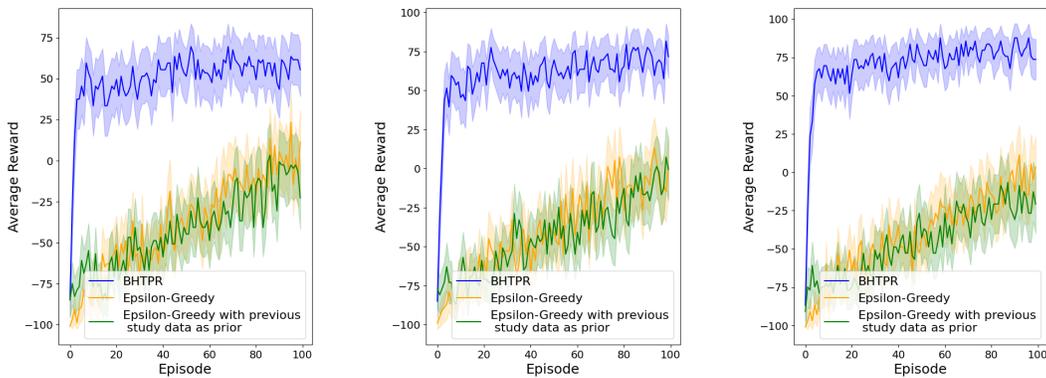


Figure 8: Comparison of initializations for algorithm parameter  $w(s)$  in the Grid MDP.

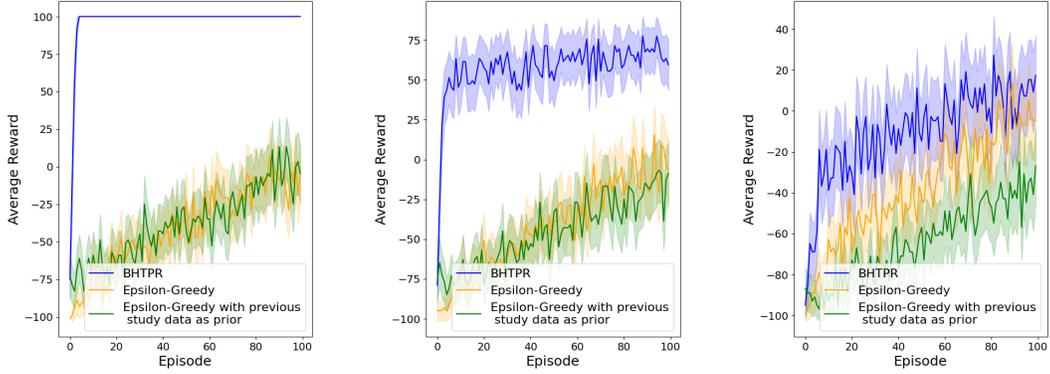
333 **E Chainworld results with alternative reward function.**

334 The plots below reflect the same set up as in the body of the paper, but with a large, negative  
 335 disengagement reward. As discussed in Sec. 4.1, a large negative reward in the disengagement state  
 336 speeds up learning, especially for the epsilon-greedy baseline.



(a) Chainworld, 50-step episodes      (b) Chainworld, 30-step episodes      (c) Chainworld, 20-step episodes

Figure 9: **Episode length.** In the chainworld, maximum episode length does not have a large effect on performance.

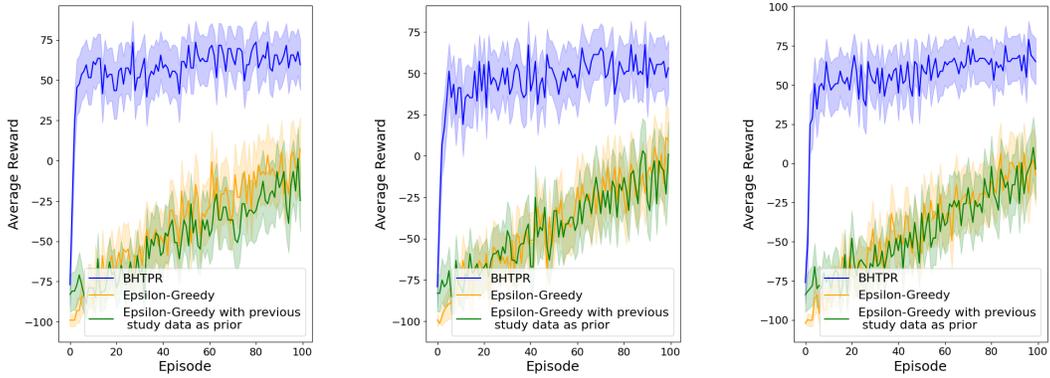


(a) Chainworld, perfect prior study data

(b) Chainworld, prior study data incorrect for states  $\geq 14$

(c) Chainworld, prior study data incorrect for states  $\geq 8$

Figure 10: **Correctness of prior study data.** Consistent with the reward function described in the body of the paper, BHTPR performs best when the prior study data matches the current environment (a), but also outperforms or matches baselines when the data does not match the environment of interest in all states.

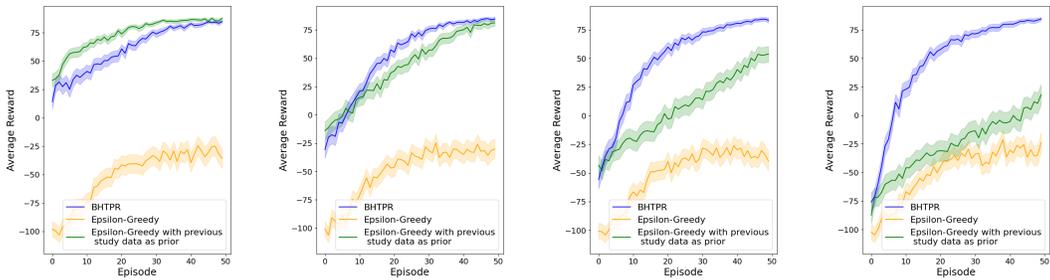


(a) Second row plot 1

(b) Second row plot 2

(c) Second row plot 3

Figure 11: **Stochasticity** Results are similar across levels of stochasticity, with BHTPR performing well in all cases.



(a) 75% correct 25% incorrect.

(b) 50% correct 50% incorrect.

(c) 25% correct 75% incorrect.

(d) 0% correct 100% incorrect.

Figure 12: **Distribution shift** When the states with corrupted data have a transition distribution that is close to correct, using the data as a prior on  $T(s, a, \cdot)$  performs best, however when the data are entirely right or wrong, our method outperforms. Plots from Grid MDP environment.