

Reinforcement learning with Demonstrations from Mismatched Task under Sparse Reward

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Reinforcement learning often suffer from the sparse reward issue in
2 real-world robotics problems. Learning from demonstration (LfD) is an effective
3 way to eliminate this problem, which leverages collected expert data to aid
4 online learning. Prior works often assume that the learning agent and the expert
5 aim to accomplish the same task, which requires collecting new data for every
6 new task. In this paper, we consider the case where the target task is mismatched
7 from but similar with that of the expert. [Such setting can be challenging and we
8 found existing LfD methods can not effectively guide learning in mismatched new
9 tasks with sparse rewards.](#) We propose conservative reward shaping from demon-
10 stration (CRSfD), which shapes the sparse rewards using estimated expert value
11 function. To accelerate learning processes, CRSfD guides the agent to conserva-
12 tively explore around demonstrations. Experimental results of robot manipulation
13 tasks show that our approach outperforms baseline LfD methods when transfer-
14 ring demonstrations collected in a single task to other different but similar tasks.

15 **Keywords:** Sparse Reward Reinforcement Learning, Learn from Demonstration,
16 Task Mismatch

17 1 Introduction

18 Reinforcement learning has been applied to various real-world tasks, including robotic manipulation
19 with large state-action spaces and sparse reward signals [1]. In these tasks, standard reinforcement
20 learning tends to perform a lot of useless exploration and easily fall into local optimal solutions.
21 To eliminate this problem, previous works often use expert demonstrations to aid online learning,
22 which adopt some successful trajectories to guide the exploration process [2, 3].

23 However, standard learning from demonstration algorithms often assume that the target learning task
24 is exactly same with the task where demonstrations are collected [4, 5, 6]. Under this assumption,
25 experts need to collect the corresponding demonstration for each new task, which can be expensive
26 and inefficient. In this paper, we consider a new learning setting where expert data is collected
27 under a single task, while the agent is required to solve different new tasks. For instance as shown
28 in Figure 1, a robot arm aims to solve peg-in-hole tasks. The demonstration is collected on a certain
29 type of hole while the target tasks have different hole shapes (changes in environmental dynamics)
30 or position shifts (changes in reward function). This can be challenging as agents cannot directly
31 imitate those demonstrations from mismatched tasks due to dynamics and reward function changes.
32 However, compared to learning from scratch, those demonstrations should still be able to provide
33 some useful information to help exploration.

34 To address the issue of learning with demonstrations from mismatched task, previous works in
35 imitation learning consider agent dynamics mismatch and rely on state-only demonstrations [7, 8, 9].
36 However, this approach has an implicit assumption that the new task share the same reward function
37 as the original task [10]. Hester et al. and Vecerik et al. [11, 3] receive sparse rewards in the

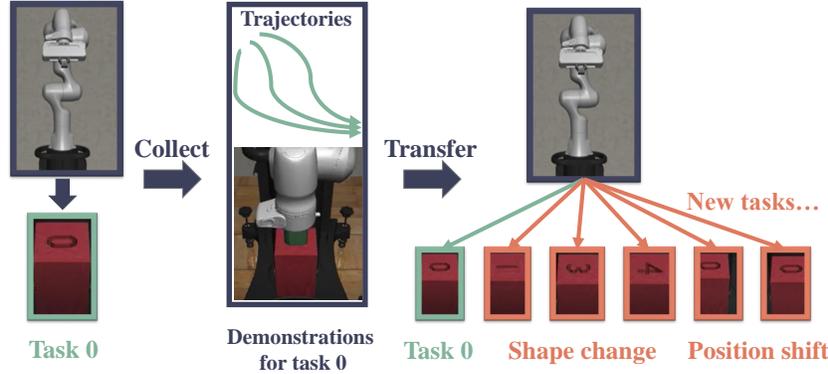


Figure 1: Illustration of our motivation. Demonstrations collected on a single original task are transferred to other similar but different tasks with either environmental dynamics changes (shape change) or reward function change (position shift), and aid the learning of these tasks.

38 environment and add demos into a prioritized replay buffer. Sparse reward signal can be backward
 39 propagated during the Bellman update and thus guide the exploration. However, this propagation
 40 flow may be blocked due to the mismatch in new tasks. Another class of work [12, 13, 14] also
 41 considers that we have expert data on multitasks and utilize meta-learning methods to obtain diverse
 42 skills, and then transfer skills to new tasks. However, such a strategy requires to collect a huge
 43 expert dataset, which is expensive and inefficient. In our setting, we are only provided with a few
 44 demonstrations collected under a single task.

45 In this paper, we propose **Conservative Reward Shaping from Demonstration (CRSfD)**, which learns
 46 policies for new tasks accelerated by demonstrations collected in a single mismatched task. We use
 47 reward shaping [15, 16] to incorporate future information into single-step rewards while keeping
 48 the optimal policy unchanged. Moreover, we explicitly deal with out-of-distribution problem to
 49 encourage agent to explore around demonstrations. Experimental results of robot manipulation
 50 tasks show that our approach outperforms baseline LfD methods when learning in new tasks with
 51 mismatched demonstrations.

52 Our contributions can be summarized as follows:

- 53 • We proposed a reward shaping scheme for reinforcement learning with demonstration from
 54 mismatched task, which use estimated value function from expert demonstrations to re-
 55 shape sparse reward in new tasks.
- 56 • Built upon such scheme, we propose the conservative reward shaping from demonstration
 57 (CRSfD) algorithm to overcome the out-of-distribution problem, we regress value function
 58 of OOD states to zero and use a larger discount factor in new tasks, which guides the agent
 59 to conservatively explore around expert data.
- 60 • We conduct simulation and real world experiments of robot insertion tasks with mis-
 61 matched demonstrations. The results show that CRSfD effectively guide the exploration
 62 process in new tasks and reach a higher sample efficiency and convergence performance.

63 2 Related Works

64 **Learning from demonstration** A prominent research subject is how to leverage expert data to as-
 65 sist reinforcement learning. Imitation learning (IL) is a broad family of such algorithms that enforce
 66 agents to directly imitate the expert. Behavior cloning (BC) is the simplest IL algorithm which
 67 greedily imitates the step-wise action of the expert and can fall into the problem of distributional
 68 shift [17]. Inverse reinforcement learning [18, 4] and adversarial imitation learning [6] infer the ex-
 69 pert’s reward function and learn the corresponding optimal policy jointly. [The above IL algorithms](#)
 70 [assume environment rewards are not available, hence their performances are upper-bounded by that](#)

71 of experts [19]. Another line of work makes use of reward feedback from environment and lever-
 72 ages expert demonstration data to overcome the sparse reward issue or learn more natural behaviors.
 73 Vecerik et al. [3] add demonstration into a prioritized replay buffer. Rajeswaran et al. [20] add
 74 a behavioral cloning loss to the policy to speed up exploration and learn more natural and robust
 75 behaviors. Chen et al. [21] use generative models on single step transition to reshape reward of the
 76 original task. However, standard learning from demonstration algorithms always requires demon-
 77 strations to be collected under the same task and act nearly optimal under this task, which is not
 78 suitable for our setting.

79 **Generalization of demonstrations** There are a few works relaxing the requirements for demonstra-
 80 tions to achieve generalization of demonstrations from different aspects. Some works assume that
 81 demonstrations are collected by a sub-optimal policy under the same task [22, 23], early work [23]
 82 requires manually ranking of trajectories and later works [24, 25] move the needs for rankings by
 83 actively adding noise to demonstrations along with automatical ranking. Cao et al. [26, 27] assume
 84 that the demonstrations are a mixture of different experts and use a classifier to separate out the
 85 more feasible expert data for the new task. Other works [10, 28, 29] assume that the target task
 86 has different agent dynamics to the task where demonstrations are collected, so they only match the
 87 state sequence of demonstrations or use an inverse dynamic model to recover the action between two
 88 states in the new task. In our work, we further consider new tasks with the environment dynamics
 89 mismatch as well as reward function mismatch. Another branch of related works are meta imitation
 90 learning algorithms, which assume that we have expert data on multitasks and utilize meta-learning
 91 methods to solve new tasks in zero-shot or few shots adaption [12, 13]. However, such a strategy
 92 usually necessitates a huge expert dataset which may be expensive and inefficient. Differently, we
 93 consider the problem where only a small number of demonstrations collected in a single task are
 94 provided, and the agent needs to use them to accelerate the learning of other similar but different
 95 tasks.

96 3 Problem Statement

97 In our problem setting, we have collected a few demonstrations under a single task and want to utilize
 98 these data in reinforcement learning for other similar but different tasks. A task can be formalized
 99 as a standard Markov decision process MDP, which is modeled as $M_i = (S, A, P_i, R_i, \gamma_i)$. The
 100 task where demonstrations are collected is denoted as $M_0 = (S, A, P_0, R_0, \gamma_0)$, and the new tasks
 101 we target to solve are denoted as $M_i = (S, A, P_i, R_i, \gamma_i), i \geq 1$. S and A are the shared state
 102 space and action space for each task. $P_i : S \times A \times S \rightarrow [0, 1]$ are state transition probability
 103 functions of each task, $R_i : S \times A \times S \rightarrow \mathbb{R}$ are reward functions for task M_i , describing the
 104 natural reward signal in each task. Due to differences of environment and agent dynamics, P_i
 105 and R_i often varied between different tasks. γ_i is discounted factor of M_i which reflects how
 106 much we care about future, typically set to a constant slightly lower than 1. A policy $\pi_i : S \rightarrow$
 107 A defines a probability distribution in action space. For a task M_i and a policy π_i , state value
 108 function $V_i^{\pi_i}(s) = \mathbb{E}_{s_0=s, \pi_i} [\sum \gamma_i^t R_t(s_t, a, s_{t+1})]$ estimates the discounted cumulative reward of the
 109 task under this policy π_i . $V_i^*(s)$ estimates the discounted cumulative rewards for state s under the
 110 optimal policy π_i .

111 As many works [30, 31] point out, directly applying RL in a sparse reward environment can be
 112 sample inefficient and fail to find a good solution. In this work, we want to make use of the demon-
 113 strations $D : (\tau_0, \tau_1, \dots)$ collected in task M_0 to facilitate reinforcement learning for the different
 114 but similar new tasks M_i . Note each trajectory τ_k contains a sequence of state action transitions
 115 $[s_0, a_0, s_1, a_1, \dots, s_t, a_t]$ in task M_0 .

116 **Challenges** There are two key issues when leveraging demonstrations from mismatched tasks. First,
 117 how to get effective guidance from these mismatched demonstrations? Although we should not
 118 purely imitate these demonstrations, we do need to obtain some useful guidance from them to ac-
 119 celeration exploration in new tasks with sparse rewards. Second, since our goal is to maximize the
 120 reward defined under the new task, guidance from mismatched demonstrations should not influence
 121 the optimality of the learned policy in new tasks.

122 4 Conservative Reward Shaping from Demonstrations

123 Provided with demonstrations in a particular task $M_0 : (S, A, P_0, R_0, \gamma_0)$, we aim to help the re-
124 inforcement learning process of different tasks $M_1, M_2 \dots M_K$, which may have different transition
125 functions $P_k(s'|s, a)$ and reward functions $R_k(s, a, s')$. In this work, we use SAC [32] as our base
126 reinforcement learning algorithm as it holds an excellent exploration mechanism which leads to
127 higher sample efficiency than policy gradient algorithms [33, 34] and is shown to perform well on
128 continuous action tasks [32]. Nevertheless, it is also possible to base our method on other RL al-
129 gorithms including on-policy ones. To make use of expert demonstrations, DDPGfD [3] proposes
130 a mechanism compatible with the off-policy method, which adds the demonstration data into the
131 replay buffer with prioritized sampling. Under such framework, the sparse reward signal can prop-
132 agate back along the expert trajectory to guide the agent. By combining SAC and DDPGfD [3], we
133 obtain the backbone of our method and labeled as SACfD, which is also our best baseline method.

134 4.1 Reward Conflict under Mismatched Task Setting

135 Although LfD methods such as SACfD benefit from demonstrations in sparse reward reinforcement
136 learning, they may not benefit from demonstrations when the target tasks are mismatched from that
137 of the expert. When following the demonstrations, agent may consistently fail and can not get any
138 sparse rewards signals. As failure time increases, agent may consider expert trajectories to have
139 low value since few rewards are received. The agent will then avoid following the expert and the
140 demonstrations cannot provide effective guidance, resulting in inefficient exploration in the whole
141 free space.

142 Although totally following the demonstrations may not be able to receive any sparse reward in new
143 tasks, it can still provide useful exploration directions since in our settings the new tasks are similar
144 to the original one. We formally introduce our method as conservative reward shaping from demon-
145 stration (CRSfD). Intuitively, CRSfD assigns appropriate reward signals along the demonstrations
146 to efficiently guide the agent towards the goal, and allows exploration around the goal to maintain
147 optimally. Details are described in the following subsection.

148 4.2 Conservative Reward Shaping from Demonstrations(CRSfD)

149 **Reward Shaping with Value Function** Reward shaping [15] provides an elegant way to modify re-
150 ward function while keeping the optimal policy unchanged. Given original MDP M and an arbitrary
151 potential function $\Phi : S \rightarrow \mathbb{R}$, we can reshape the reward function to be:

$$R'(s, a, s') = R(s, a, s') + \gamma\Phi(s') - \Phi(s), s' \sim P(\cdot|s, a) \quad (1)$$

152 Denote the new MDP as $M' = (S, A, P, R', \gamma)$ obtained by replacing reward function R in M to R' .
153 Ng et al. [15] proved that the optimal policy $\pi_{M'}$ on M' and the optimal policy π_M^* on the original
154 MDP M_0 are the same: $\pi_{M'}^* = \pi_M^*$. Furthermore, the optimal state-action function $Q_{M'}^*(s, a)$ and
155 value function $V_{M'}^*(s)$ are shifted by $\Phi(s)$:

$$Q_{M'}^*(s, a) = Q_M^*(s, a) - \Phi(s), \quad V_{M'}^*(s) = V_M^*(s) - \Phi(s) \quad (2)$$

156 In particular, Ng et al. [15] pointed out that when the potential function is chosen as the optimal
157 value function of the original MDP $\Phi(s) = V_M^*(s)$, then the new MDP M' becomes trivial to solve.
158 What remained for the agent is to choose each time-step's action greedily, because the transformed
159 single-step reward already contains all the long-term information for decision making.

160 **Conservative Value Function Estimation** The reward shaping method provides a principled way
161 to guide the agent with useful future information and keep the optimal policy unchanged. Ideally,
162 an accurate $\Phi_i(s) = V_{M_i}^*(s)$ will lead to simple and optimal policy in new MDP M' , but a perfect
163 $\Phi_i(s) = V_{M_i}^*(s)$ is unavailable in advance. Practically, we estimate a $\tilde{V}_{M_0}^D \approx V_{M_0}^*(s)$ using demon-
164 strations from task M_0 by Monte-Carlo regression and treat $\tilde{V}_{M_0}^D(s)$ as a prior guess of $V_{M_i}^*(s)$. We
165 then shape the sparse reward in the new task M_i to:

$$R'_i(s, a, s') = R_i(s, a, s') + \gamma\tilde{V}_{M_0}^D(s') - \tilde{V}_{M_0}^D(s) \quad (3)$$

166 However, demonstration trajectories only cover a small part of the state space. For out-of-
167 distribution states, estimated $\tilde{V}_{M_0}^D$ may output random values and lead to random single-step re-
168 ward after reward shaping, which may mislead the agent. We make two improvements over the

169 above reward shaping method to encourage the agent to explore around the demonstrations conser-
 170 vatively: (1) Regress value function $\tilde{V}_{M_0}^*(s)$ of the out-of-distribution states to 0, thus discouraging
 171 exploration far from demonstrations. The OOD states are sampled randomly from free space. (2) In-
 172 creasing the discount factor γ_i in new tasks. From equation 3, we can find that increasing γ_i will give
 173 higher single-step reward for state with large $V_\theta(s')$ in the original task, thus encourages exploration
 174 around demonstrations. Our method can be summarized as follows: (D stands for demonstration
 175 buffer, S stands for free space, $\gamma_i > \gamma_0$):

Algorithm 1 Conservative Value Function Estimation

Input: Demonstration transitions, demo discount factor γ_0 , new task discount factor $\gamma_k(\gamma_k > \gamma_0)$, regression steps n_r , scale factor λ .

Initialization: Initialize value function $V_\theta(s)$

Monte-Carlo policy evaluation on demonstrations, Calculate cumulative reward for states in demos using γ_0 : $V_{M_0}^D(s) = \sum_{i=t}^T \gamma_0^{i-t} r_i$

for n in regression steps n_r **do**

 Sample minibatch B_1 from demo buffer D with regression target $V_{M_0}^D(s) = \sum_{i=t}^T \gamma_0^{i-t} r_i$. Sample minibatch B_2 from whole free space S with regression target 0.

 perform regression: $\theta = \arg \min_{\theta} \left[\mathbb{E}_{s_t \sim B_1} (V_\theta(s_t) - \sum_{i=t}^T \gamma_0^{i-t} r_i)^2 + \lambda \mathbb{E}_{s_t \sim B_2} (V_\theta(s_t) - 0)^2 \right]$

end for

Shaping reward with γ_k : $R'_i(s, a, s') = R_i(s, a, s') + \gamma_k V_\theta(s') - V_\theta(s)$.

Perform SACfD update. (details can be found in appendix.)

176 **Conservative Properties** In the last paragraph, we introduced some conservative techniques and
 177 give some intuitively explanations why those improvements can encourage exploration around
 178 demonstrations under the proposed reward shaping framework. The following theorem can quantize
 179 the benefits of proposed methods.

180 **Theorem 1** For task M_0 with transition T_0 and new task M_k with transition T_k , define to-
 181 tal variation divergence $D_{TV}(s, a) = \sum_{s'} |T_0(s'|s, a) - T_k(s'|s, a)| = \delta$. If we have $\delta <$
 182 $(\gamma_k - \gamma_0) \mathbb{E}_{T_2(s'|s, a)} [V_{M_0}^D(s')] / \gamma_0 \max_{s'} V_{M_0}^D(s')$, then following the expert policy in new task will
 183 result in immediate reward greater then 0:

$$\mathbb{E}_{a \sim \pi(\cdot|s)} r'(s, a) \geq (\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s)} [V_{M_0}^D(s')] - \gamma_0 \delta \max_{s'} V_{M_0}^D(s') > 0 \quad (4)$$

184 Detailed proof can be found in Appendix 7.5. The above theorem indicates that for similar but
 185 different tasks (δ smaller than the threshold), exploration along demonstrations will lead to positive
 186 immediate rewards which guide the learning process.

187 **Conservative Reward Shaping from Demonstrations** After reward shaping by demonstrations
 188 from mismatched task, we perform online learning based on SACfD as described in Section 4. Pseu-
 189 docode can be found in supplementary materials. Although the estimated $\tilde{V}_{M_0}^D(s)$ can be inaccurate,
 190 it still provides enough future information, thus facilitates exploration for the agent. Moreover, nice
 191 theoretical properties of reward shaping guarantees that we will not introduce bias to the learned
 192 policy in new tasks.

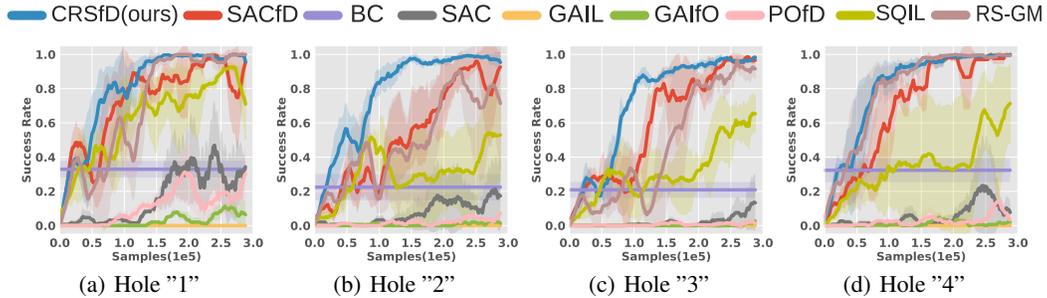


Figure 2: Evaluation of algorithms on 4 new tasks with demonstrations from task "0". The solid line corresponds to the mean of success rate over 5 random seeds and the shaded region corresponds to the standard deviation. Y-axis reflects success rate range in $[0, 1]$, X-axis reflects interaction steps range in $[0, 3e5]$.

193 5 Experimental Results

194 We perform experimental evaluations of the proposed CRSfD method and try to answer the follow-
195 ing two questions: Can CRSfD help the exploration of similar sparse-rewarded tasks with demon-
196 strations from a mismatched task? Will CRSfD introduce bias to the learned policy in new tasks?

197 We choose the robot insertion tasks for our experiments, which has natural sparse reward signals:
198 successfully inserting peg into hole get a reward +1, otherwise 0. We perform both simulation and
199 real world experiments. The simulation environment is built under robosuite framework [35] pow-
200 ered by Mujoco physics simulator [36]. We construct a series of similar tasks where the holes have
201 different shapes and unknown position shifts, reflecting changes in dynamics and reward functions
202 respectively, as shown in Figure 1. Then we verify the effectiveness of CRSfD under the following 2
203 settings: (1) Transfer collected demonstrations to similar insertion tasks with environment dynamics
204 mismatch. (2) Transfer collected demonstrations to similar tasks with both environment dynamics
205 and reward function mismatch. Finally, we address the sim-to-real issue and deploy the learned pol-
206 icy on a real robot arm to perform insertion tasks with various shapes of holes in the real world. We
207 use Franka Panda robot arm in both simulation and real world. The comparison baseline algorithms
208 are chosen as follows:

- 209 • **Behavior Cloning [17]**: Just ignore the task mismatch and directly perform behavior
210 cloning of the demonstrations.
- 211 • **SAC [32]**: A SOTA standard RL method which does not use the demonstrations and di-
212 rectly learn from scratch in the target tasks.
- 213 • **GAIL [6]**: Use adversarial training to recover the policy that generates the demonstrations,
214 which alleviates the distributional shift problem of behavior cloning.
- 215 • **GAIfO [8]**: A variant of GAIL which trains a discriminator with state transitions (s, s')
216 instead of (s, a) in GAIL to alleviate dynamics mismatch.
- 217 • **POfD [37]**: A variant of GAIL which combines the intrinsic reward from discriminator
218 and extrinsic reward from the new task.
- 219 • **SQIL [38]**: An effective off-policy imitation learning algorithm that adds demonstrations
220 with reward +1 to the buffer and assign reward 0 to all agent experiences.
- 221 • **SACfD [32, 3]**: Incorporate effective demonstration replay mechanism from [3] with SAC
222 as described in Section 4, which is also the best baseline as well as the backbone of our
223 method.
- 224 • **RS-GM [21]**: Reward Shaping using Generative Models, which is an extension on discrete
225 reward shaping methods [39, 40]. After learning a discriminator $D_\phi(s, a)$, they shape the
226 reward into $R'(s, a) = R(s, a) + \gamma\lambda D_\phi(s', a') - \lambda D_\phi(s, a)$.

227 5.1 Simulation experiments

228 We set a nominal hole position as the original-point of our Cartesian coordination. Observable states
229 include robot proprioceptive information such as joint and end-effector position and velocity. Action
230 space includes the 6d pose change of the robot end-effector in 10 Hz, followed with a Cartesian
231 impedance PD controller running at a high frequency. Only a sparse reward +1 is provided when
232 the peg is totally inserted inside the hole. Demonstrations are collected by a sub-optimal RL policy
233 trained with SAC in task M_0 under carefully designed dense reward, where the hole has shape "0".
234 This process can be replaced by manual collection in the real world. Then demonstrations are tagged
235 with the corresponding sparse reward. We collected 40 demonstration trails with 50 time steps each.

236 **Setting 1: Tasks with environment dynamics mismatch.** To reflect environment dynamics
237 changes of the tasks, we create experiments domain on insertion tasks with holes of various shapes
238 in the simulator, represented by different digit numbers, as shown in Figure 1. Different shapes
239 of holes will encounter different contact mode thus lead to different environmental dynamics. We
240

241 collect demonstrations from hole "0", and our method use them to help training similar new tasks
 242 with various hole shapes from digit 1 to 4.

243 **Analysis 1:** The comparison results of CRSfD and baseline algorithms under the above setting are
 244 show in Figure 2 . As we expected, the simplest BC algorithm simply imitates the expert action
 245 of the original task and can only complete the insertion with a small chance. The SAC algorithm
 246 does not make use of the demonstration data and conducts a lot of useless exploration, which leads
 247 to poor performance. GAIL algorithm and its variants GAIfo, POfo also fail for most of times as
 248 they try to purely imitate the demonstration collected in the mismatched task. SQIL ignores the
 249 reward in the new task and only obtains a limited success rate. **SACfD can not be effectively guided**
 250 **by demonstrations from the mismatched task under sparse reward.** Our proposed CRSfD provide
 251 guidance through reward shaping, and consistently achieves the best performance on all the four
 252 insertion tasks with different hole shapes.

253 **Setting 2: Tasks with both dynamics and reward function mismatch** Next, we consider more
 254 challenging scenarios where we aim to transfer the demonstrations to new tasks with both environ-
 255 mental dynamics mismatch and reward function mismatch. We assume that the hole has unknown
 256 random shifts relative to the nominal position, thus the reward function changes. At the beginning
 257 of each episode, the hole is uniformly initialized in a square area centered at the nominal position.
 258 This can be challenging because the robot is 'blind' to these unknown offsets and requires further
 259 search for the entrance of the hole. Practically, we collected demonstrations from task with hole "0"
 260 with fixed hole position, and transfer to new tasks with random hole shifts and different hole shapes.

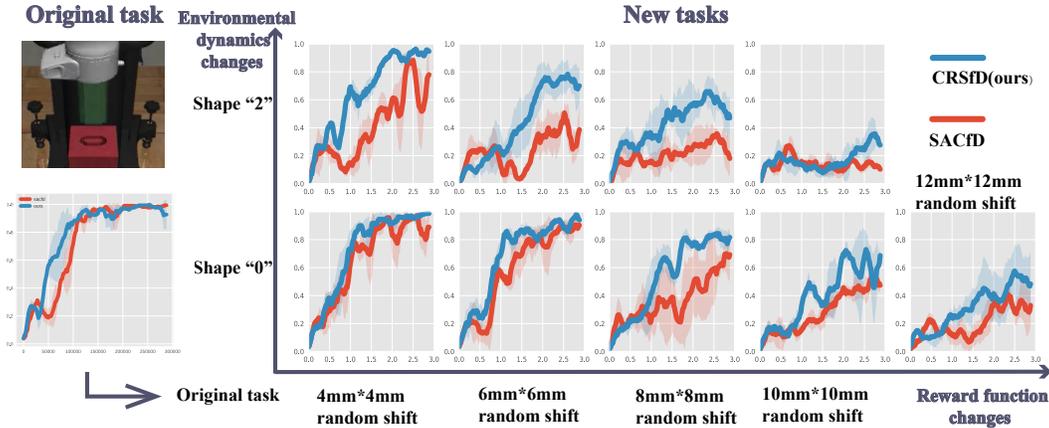


Figure 3: Evaluations of CRSfD and the best baseline SACfD. The solid line corresponds to the mean of success rate over 3 random seeds and the shaded region corresponds to the standard deviation. X-coordinate reflects changes in reward functions and Y-coordinate reflects changes in environmental dynamics. Our algorithm outperforms baseline with increasing margins as the task changes become larger.

261 **Analysis 2:** We compare our algorithm with the best baseline algorithm SACfD under varying de-
 262 grees of environmental dynamics and reward function changes, as shown in the Figure 3. Due to
 263 space limit, more comparison can be found in Figure 7 in appendix. The x-coordinate represents
 264 the increasing changes of the reward function, where the random range of the holes becomes larger
 265 (from 4mm*4mm, 6mm*6mm, to 8mm*8mm). The y-coordinate represents increasing environmen-
 266 tal dynamics change, from hole "0" in its original shape to hole "2" in a different shape. Straight-
 267 forwardly, coordinate origin can represent the original task where demonstrations are collected, and
 268 a 2d coordinates $[x, y]$ represents a new task with varying degree of mismatch.

269 From Figure 3, we can observe that when applying to the original task or very similar task such as
 270 [4mm, shape '0'], our method has a similar performance to the SACfD baseline. When the task
 271 changes become greater (e.g, [8mm, shape '0'], [4mm, shape '2'], [6mm, shape '2'], [8mm, shape
 272 '2']), **SACfD gradually lose the guidance from original demonstrations as task mismatched more**
 273 **significantly**, while CRSfD achieves significant performance gains with help of the conservative
 274 reward shaping using estimated value function.

275 **Ablation study** As mentioned in section 4.2, we make two improve-
 276 ments over the reward shaping method to encourage the agent to explore
 277 around the demonstrations conservatively. (1) Regress value function of
 278 OOD states to zero. (2) Use a larger discount factor in new tasks. We
 279 ablate these 2 improvements and compare their performance. Ablations
 280 are tested under new task with hole shape “3”, results for other shapes
 281 can be found in the supplementary materials. As shown in Figure 5.1,
 282 compared to original CRSfD algorithm, moving away either of these 2
 283 techniques will lead to a performance drop, where the agent needs to
 284 take more effort in exploration.

285 5.2 Real World Experiments

286 After completing the insertion tasks of various-shaped holes in the sim-
 287 ulator, we deploy the policy to the real robotic arm. To overcome the
 288 sim-to-real problem, we use domain randomization in the simulation.
 289 The initial position of the robot arm end-effector and holes are randomized in a 6cm*6cm*6cm
 290 space and 2mm*2mm plane respectively, and the friction coefficient of the object is also random-
 291 ized in [1, 2]. We use a real Franka Panda robot arm and 3d print the holes corresponding to digit
 292 numbers “0-4”. Holes are roughly in sizes of 4cm*4cm, with a 1mm clearance between the peg and
 293 the hole. We performed 25 insertion trials under each shape of hole, and counted their success rates
 294 separately, as shown in the table 1. The robot achieves high success rate in all tasks.

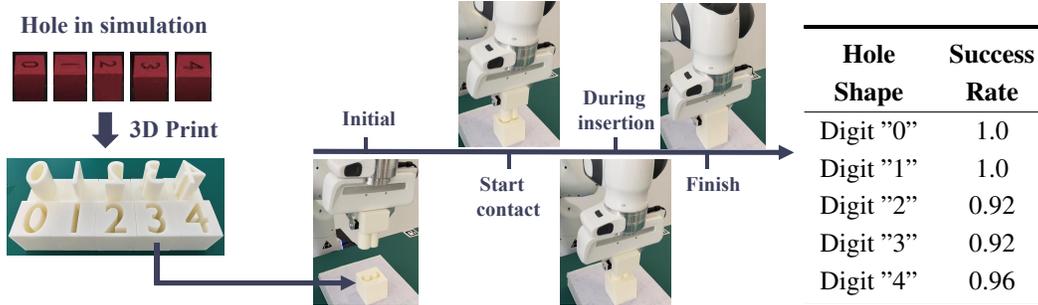


Figure 5: Real world robot insertion experiments.

295 6 Conclusion

296 **Summary.** In this paper, we studied the problem of reinforcement learning with demonstrations
 297 from mismatched tasks under sparse rewards. Our key insight is that, although we should not purely
 298 imitate the mismatched demonstrations, we can still get useful guidance from the demonstrations
 299 collected in a similar task. Concretely, we proposed conservative reward shaping from demonstra-
 300 tions (CRSfD) which uses reward shaping by estimated value function of a mismatched expert to
 301 incorporate useful future information to augment the sparse reward, with conservativeness tech-
 302 niques to handle out-of-distribution issues. Simulation and real world robot insertion experiments
 303 show the effective of proposed method under tasks varied in environmental dynamics and reward
 304 functions.

305 **Limitations and Future works.** Provided with demonstrations from a mismatched task, our pro-
 306 posed method aids the online learning process for each new task separately. However, one may need
 307 to learn a policy to solve multiple new tasks at the same time, and exploration in these tasks may
 308 benefit each other. So future works include using demonstrations to accelerate the joint learning pro-
 309 cess of multiple tasks. Another limitation is that our method is only applicable to new tasks similar
 310 to original task. The effectiveness of CRSfD gradually decays when the tasks differ too much from
 311 the original task so that the demonstrations do not contain any useful information. It also worth to
 312 mention that the whole algorithm pipeline should be able to be implemented directly on hardware,
 313 which is a promising research direction.

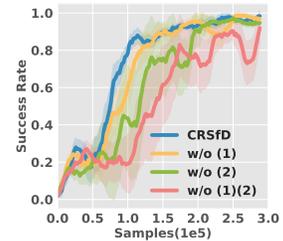


Figure 4: Ablation studies of the conservativeness techniques. (1) means regressing value function to zero for OOD states. (2) means setting larger discount factors.

References

- 314
- 315 [1] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki,
316 A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation.
317 *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- 318 [2] A. Nair, B. McGrew, M. Andrychowicz, W. Zaremba, and P. Abbeel. Overcoming exploration
319 in reinforcement learning with demonstrations. In *2018 IEEE international conference on*
320 *robotics and automation (ICRA)*, pages 6292–6299. IEEE, 2018.
- 321 [3] M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe,
322 and M. Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics
323 problems with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017.
- 324 [4] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, et al. Maximum entropy inverse reinforce-
325 ment learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- 326 [5] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Pro-*
327 *ceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- 328 [6] J. Ho and S. Ermon. Generative adversarial imitation learning. *Advances in neural information*
329 *processing systems*, 29, 2016.
- 330 [7] Y. Schroecker and C. L. Isbell. State aware imitation learning. *Advances in Neural Information*
331 *Processing Systems*, 30, 2017.
- 332 [8] F. Torabi, G. Warnell, and P. Stone. Generative adversarial imitation from observation. *arXiv*
333 *preprint arXiv:1807.06158*, 2018.
- 334 [9] W. Sun, A. Vemula, B. Boots, and D. Bagnell. Provably efficient imitation learning from
335 observation alone. In *International conference on machine learning*, pages 6036–6045. PMLR,
336 2019.
- 337 [10] T. Gangwani and J. Peng. State-only imitation with transition dynamics mismatch. *arXiv*
338 *preprint arXiv:2002.11879*, 2020.
- 339 [11] T. Hester, M. Vecerik, O. Pietquin, M. Lanctot, T. Schaul, B. Piot, D. Horgan, J. Quan,
340 A. Sendonaris, I. Osband, et al. Deep q-learning from demonstrations. In *Proceedings of*
341 *the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 342 [12] V. H. Pong, A. Nair, L. Smith, C. Huang, and S. Levine. Offline meta-reinforcement learning
343 with online self-supervision. *arXiv preprint arXiv:2107.03974*, 2021.
- 344 [13] T. Z. Zhao, J. Luo, O. Sushkov, R. Pevcevičute, N. Heess, J. Scholz, S. Schaal, and S. Levine.
345 Offline meta-reinforcement learning for industrial insertion. *arXiv preprint arXiv:2110.04276*,
346 2021.
- 347 [14] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A
348 benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
349 *Robot Learning*, pages 1094–1100. PMLR, 2020.
- 350 [15] A. Y. Ng, D. Harada, and S. Russell. Policy invariance under reward transformations: Theory
351 and application to reward shaping. In *Icml*, volume 99, pages 278–287, 1999.
- 352 [16] C.-A. Cheng, A. Kolobov, and A. Swaminathan. Heuristic-guided reinforcement learning.
353 *Advances in Neural Information Processing Systems*, 34, 2021.
- 354 [17] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured predic-
355 tion to no-regret online learning. In *Proceedings of the fourteenth international conference*
356 *on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Pro-
357 ceedings, 2011.

- 358 [18] A. Y. Ng, S. J. Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1,
359 page 2, 2000.
- 360 [19] D. Rengarajan, G. Vaidya, A. Sarvesh, D. Kalathil, and S. Shakkottai. Reinforcement
361 learning with sparse rewards using guidance from offline demonstration. *arXiv preprint*
362 *arXiv:2202.04628*, 2022.
- 363 [20] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine.
364 Learning complex dexterous manipulation with deep reinforcement learning and demonstra-
365 tions. *arXiv preprint arXiv:1709.10087*, 2017.
- 366 [21] Y. Wu, M. Mozifian, and F. Shkurti. Shaping rewards for reinforcement learning with imperfect
367 demonstrations using generative models. In *2021 IEEE International Conference on Robotics*
368 *and Automation (ICRA)*, pages 6628–6634. IEEE, 2021.
- 369 [22] Y. Gao, H. Xu, J. Lin, F. Yu, S. Levine, and T. Darrell. Reinforcement learning from imperfect
370 demonstrations. *arXiv preprint arXiv:1802.05313*, 2018.
- 371 [23] D. Brown, W. Goo, P. Nagarajan, and S. Niekum. Extrapolating beyond suboptimal demon-
372 strations via inverse reinforcement learning from observations. In *International conference on*
373 *machine learning*, pages 783–792. PMLR, 2019.
- 374 [24] L. Chen, R. Paleja, and M. Gombolay. Learning from suboptimal demonstration via self-
375 supervised reward regression. *arXiv preprint arXiv:2010.11723*, 2020.
- 376 [25] D. S. Brown, W. Goo, and S. Niekum. Better-than-demonstrator imitation learning via
377 automatically-ranked demonstrations. In *Conference on robot learning*, pages 330–359.
378 PMLR, 2020.
- 379 [26] Z. Cao and D. Sadigh. Learning from imperfect demonstrations from agents with varying
380 dynamics. *IEEE Robotics and Automation Letters*, 6(3):5231–5238, 2021.
- 381 [27] Z. Cao, Y. Hao, M. Li, and D. Sadigh. Learning feasibility to imitate demonstrators with
382 different dynamics. *arXiv preprint arXiv:2110.15142*, 2021.
- 383 [28] I. Radosavovic, X. Wang, L. Pinto, and J. Malik. State-only imitation learning for dexterous
384 manipulation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems*
385 *(IROS)*, pages 7865–7871. IEEE, 2020.
- 386 [29] F. Liu, Z. Ling, T. Mu, and H. Su. State alignment-based imitation learning. *arXiv preprint*
387 *arXiv:1911.10947*, 2019.
- 388 [30] J. Oh, Y. Guo, S. Singh, and H. Lee. Self-imitation learning. In *International Conference on*
389 *Machine Learning*, pages 3878–3887. PMLR, 2018.
- 390 [31] M. Riedmiller, R. Hafner, T. Lampe, M. Neunert, J. Degraeve, T. Wiele, V. Mnih, N. Heess,
391 and J. T. Springenberg. Learning by playing solving sparse reward tasks from scratch. In
392 *International conference on machine learning*, pages 4344–4353. PMLR, 2018.
- 393 [32] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
394 deep reinforcement learning with a stochastic actor. In *International conference on machine*
395 *learning*, pages 1861–1870. PMLR, 2018.
- 396 [33] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization.
397 In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- 398 [34] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
399 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- 400 [35] Y. Zhu, J. Wong, A. Mandlekar, and R. Martín-Martín. robosuite: A modular simulation
401 framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020.
- 402 [36] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
403 *IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE,
404 2012.
- 405 [37] B. Kang, Z. Jie, and J. Feng. Policy optimization with demonstrations. In *International con-*
406 *ference on machine learning*, pages 2469–2478. PMLR, 2018.
- 407 [38] S. Reddy, A. D. Dragan, and S. Levine. Sqil: Imitation learning via reinforcement learning
408 with sparse rewards. *arXiv preprint arXiv:1905.11108*, 2019.
- 409 [39] T. Brys, A. Harutyunyan, H. B. Suay, S. Chernova, M. E. Taylor, and A. Nowé. Reinforcement
410 learning from demonstration through shaping. In *Twenty-fourth international joint conference*
411 *on artificial intelligence*, 2015.
- 412 [40] A. Hussein, E. Elyan, M. M. Gaber, and C. Jayne. Deep reward shaping from demonstrations.
413 In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 510–517. IEEE,
414 2017.

415 7 Appendix

416 7.1 Algorithm

Algorithm 2 CRSfD

Input: Env Environment for the new task M_i ; θ^π initial policy parameters; θ^Q initial action-value function parameters; $\theta^{Q'}$ initial target action-value function parameters; N target network update frequency.

Input: B^E replay buffer initialized with demonstrations. B replay buffer initialized empty. K number of pre-training gradient updates. d expert buffer sample ratio. $batch$ mini batch size.

Input: θ^V initial value function (potential function), original task discount factor γ_0 .

Output: $Q_\theta(s, a)$ action-value function (critic) and $\pi(\cdot|s)$ the policy (actor).

Estimate value function from demonstration.

for step t **in** $\{0, 1, 2, \dots, T\}$ **do**

 Sample with $batch$ transitions from B^E , calculate their Monte-Carlo return with discount factor γ_0 .

 Estimate $V_\theta(s)$ conservatively by equation ??

end for

Interact with Env .

for episode e **in** $\{0, 1, 2, \dots, M\}$ **do**

 Initialize state $s_0 \sim Env$

for step t **in** episode length $\{0, 1, 2, \dots, T\}$ **do**

 Sample action from $\pi(\cdot|s_t)$

 Get next state and natural sparse reward s_{t+1}, r_t

 Shape reward by: $r'_t = r_t + \gamma_i V(s_{t+1}, \theta^V) - V(s_t, \theta^V)$

 Add single step transition $(s_t, a_t, r'_t, s_{t+1})$ to replay buffer B .

end for

for update step l **in** $\{0, 1, 2, \dots, L\}$ **do**

 Sample with prioritization: $d * batch$ transitions from B^E , $(1 - d) * batch$ transitions from B . Concatenate them into a single batch.

 Perform SAC update for actor and critic: $L_{Actor}(\theta^\pi), L_{Critic}(\theta^Q)$.

if step $l \equiv 0 \pmod{N}$ **then**

 Update target critic using moving average: $\theta^{Q'} = (1 - \tau)\theta^{Q'} + \tau\theta^Q$

 Decrease expert buffer sample ratio: $d = d - \delta$ if $d > 0$.

end if

end for

end for

417 7.2 Implementation Details

418 We implemented our CRSfD algorithm and the baseline algorithms in PyTorch and the implemen-
419 tation can be found in the supplementary materials. Simulated environments are based on ro-
420 bosuite framework <https://github.com/ARISE-Initiative/robosuite>. Our CRSfD algo-
421 rithm is based on https://github.com/denisyarats/pytorch_sac_ae while baseline algo-
422 rithms are based on <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail> and
423 <https://github.com/ku2482/gail-airl-ppo.pytorch>.

424 7.3 Videos

425 Videos for simulated environments and real world environments can be found in the supplementary
426 materials.

427 7.4 Ablations

428 As mentioned in section 5.1, we make two improvements over the reward shaping method to en-
 429 courage the agent to explore around the demonstrations conservatively. (1) Regress value function
 430 of OOD states to zero. (2) Use a larger discount factor in new tasks.

431 We ablate these 2 improvements and compare their performance on more environments, as show in
 432 Figure 6.

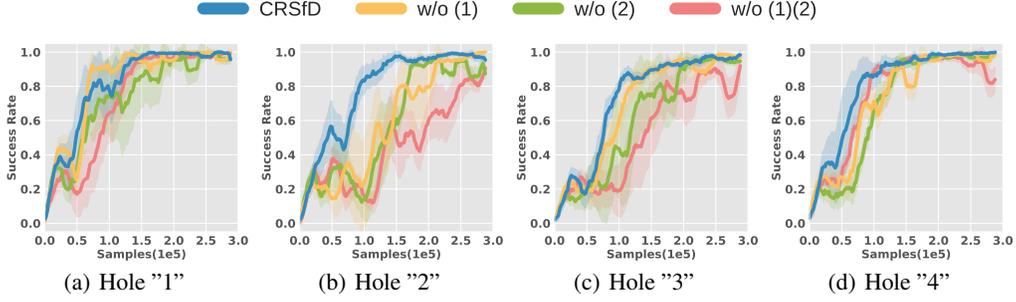


Figure 6: Ablation studies of the conservativeness techniques. (1) means regressing value function to zero for OOD states. (2) means setting larger discount factors.

433 7.5 Proof for theorem

434 **Theorem 1** For task M_0 with transition T_0 and new task M_k with transition T_k , define to-
 435 tal variation divergence $D_{TV}(s, a) = \sum_{s'} |T_0(s'|s, a) - T_k(s'|s, a)| = \delta$. If we have $\delta <$
 436 $(\gamma_k - \gamma_0) \mathbb{E}_{T_2(s'|s, a)} [V_{M_0}^D(s')]/\gamma_0 \max_{s'} V_{M_0}^D(s')$, then following the expert policy in new task will
 437 result in immediate reward greater then 0:

$$\mathbb{E}_{a \sim \pi(\cdot|s)} r'(s, a) \geq (\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s)} [V_{M_0}^D(s')] - \gamma_0 \delta \max_{s'} V_{M_0}^D(s') > 0 \quad (5)$$

438 **Proof:** For simplify, denote demonstration state value function in original task $V_{M_0}^D = V_1(s)$. Start
 439 from the reward shaping equations, and extend $V_1(s)$ for one more time step:

$$\begin{aligned} r'(s, a, s') &= r(s, a, s') + \gamma_k V_1(s') - V_1(s) \\ r'(s, a) &= r(s, a) + \gamma_k \mathbb{E}_{T_k(s'|s, a)} [V_1(s')] - V_1(s) \\ &= (\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s, a)} [V_1(s')] + (r(s, a) + \gamma_0 \mathbb{E}_{T_k(s'|s, a)} [V_1(s')] - V_1(s)) \\ &\geq (\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s, a)} [V_1(s')] + (Q^{\pi_1}(s, a) - V_1(s)) - \gamma_0 \delta \max_{s'} V_1(s') \end{aligned} \quad (6)$$

440 Take expectation on demonstration policies:

$$\mathbb{E}_{a \sim \pi(\cdot|s)} r'(s, a) \geq \mathbb{E}_{a \sim \pi(\cdot|s)} [(\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s, a)} [V_1(s')]] - \gamma_0 \delta \max_{s'} V_1(s') \quad (7)$$

441 For a sparse reward environment, we have $r(s, a) = 0$ almost everywhere:

$$\begin{aligned} \mathbb{E}_{a \sim \pi(\cdot|s)} r'(s, a) &\geq \mathbb{E}_{a \sim \pi(\cdot|s)} [(\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s, a)} [V_1(s')]] - \gamma_0 \delta \max_{s'} V_1(s') \\ &= (\gamma_k - \gamma_0) \mathbb{E}_{T_k(s'|s)} [V_1(s')] - \gamma_0 \delta \max_{s'} V_1(s') \end{aligned} \quad (8)$$

442 **7.6 Increasingly Larger Task Mismatch**

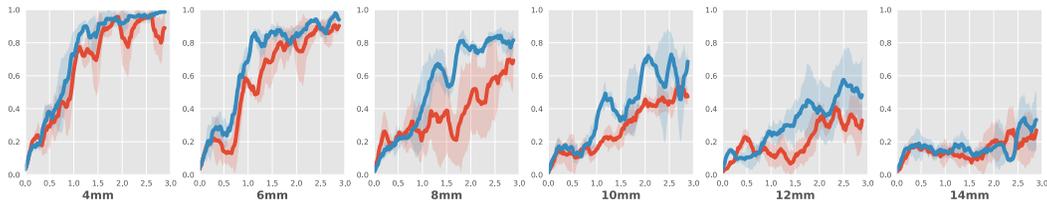


Figure 7: Increasingly larger task mismatch. Experiments are done on hole shape 0 with increasing random hole position.

443 We can observe that as task difference increases, our method first gradually outperforms baseline
444 methods. When task mismatch are too large, our method gradually loss some performance and has
445 similar performance with baselines.