
You Only Live Once: Single-Life Reinforcement Learning via Learned Reward Shaping

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement learning algorithms are typically designed to learn a performant
2 policy that can repeatedly and autonomously complete a task, typically starting
3 from scratch. However, many real-world situations operate under a different set of
4 assumptions: the goal might not be to learn a policy that can do the task repeatedly,
5 but simply to perform a new task successfully once, ideally as quickly as possible,
6 and while leveraging some prior knowledge or experience. For example, imagine
7 a robot that is exploring another planet, where it cannot get help or supervision
8 from humans. If it needs to navigate to a crater that it has never seen before in
9 search of water, it does not really need to acquire a policy for reaching craters
10 reliably, it only needs to reach this particular crater once. It must do so without the
11 benefit of episodic resets and tackle a new, unknown terrain, but it can leverage
12 prior experience it acquired on Earth. We formalize this problem setting, which
13 we call *single-life reinforcement learning* (SLRL), where an agent must complete
14 a task once while contending with some form of novelty in a single trial without
15 interventions, given some prior data. In this setting, we find that algorithms
16 designed for standard episodic reinforcement learning can struggle, as they have
17 trouble recovering from novel states especially when informative rewards are
18 not provided. Motivated by this observation, we also propose an algorithm, *Q*-
19 weighted adversarial learning (QWALE), that addresses the dearth of supervision
20 by employing a distribution matching strategy that leverages the agent’s prior
21 experience as guidance in novel situations. Our experiments on several single-
22 life continuous control problems indicate that methods based on our distribution
23 matching formulation are 20-60% more successful because they can more quickly
24 recover from novel, out-of-distribution states.

25 1 Introduction

26 When building autonomous agents for the natural world, often the goal is not to learn a performant
27 policy but rather to get something done, perhaps even suboptimally. For example, an agent exploring
28 on Mars looking for water will only need to complete its mission a single time. As another example,
29 a rescue robot will need to recover valuables from a particular burning building only once. While
30 the agent may have access to prior data about its task, a challenge arises from the fact the agent is
31 inevitably going to have to contend with some form of novelty. In the prior examples, the agent
32 on Mars may have to contend with unknown terrain and environmental conditions, and the rescue
33 agent may find certain paths in the building unpassable due to the fire. We formalize this setting as
34 the *single-life reinforcement learning* (SLRL) setting, where the agent is evaluated on its ability to
35 complete a task in a single trial autonomously without episodic resets. Importantly, the given online
36 task contains an aspect of novelty not present in the prior data although the task objective remains the
37 same. The agent’s objective is to complete the given task as quickly as possible, rather than learn a
38 policy that can repeatedly complete the task.

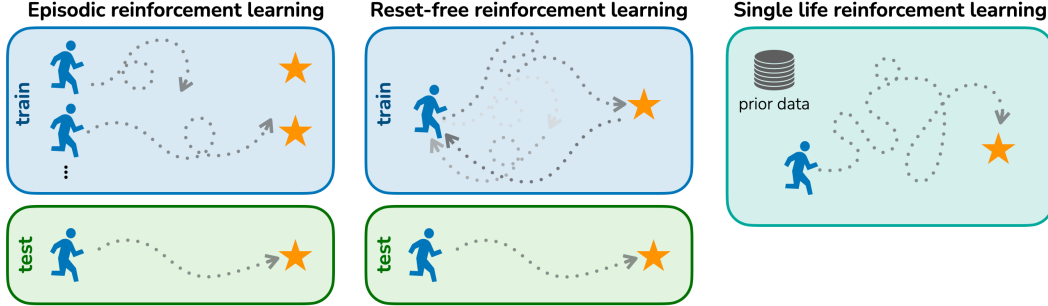


Figure 1: We study the single-life reinforcement learning (SLRL) problem, where given prior data, an agent must complete a task autonomously in a single trial in a domain with a novel distribution shift.

39 We find that algorithms designed for episodic policy learning can struggle to complete single-life
 40 tasks, even when initialized with prior data. These algorithms empirically struggle to recover from
 41 novel states. In episodic RL, the agent can rely on a reset to recover from an unfamiliar state. In
 42 contrast, in SLRL, the agent will inevitably fall off the distribution of prior data and must find its way
 43 back to a good state distribution on its own. We find that fine-tuning a pre-trained value function via
 44 online RL will not explicitly encourage the agent to get back on distribution. We hypothesize that
 45 biasing exploration towards the known distribution represented in prior data and incentivizing the
 46 agent to stay there may be suboptimal from a policy learning perspective but may enable the agent to
 47 get the task done more quickly, which is what we care about in SLRL. While shaped rewards may
 48 help the agent find its way back and ultimately complete the desired task, the agent may often find
 49 itself in a sparse reward environment or with access to rewards that are not informative enough to be
 50 guided towards task completion.

51 Adversarial imitation learning (AIL) approaches such as GAIL ([21]) can potentially provide the
 52 desired reward shaping via distribution matching. However, using existing AIL methods naively may
 53 not give the intended behavior in the SLRL setting due to two main shortcomings. First, such methods
 54 assume expert demonstrations are given as prior data, but in SLRL, we may be given suboptimal
 55 offline prior data. Second, AIL methods train the agent to match the entire distribution of prior
 56 data, which may be key to learning an optimal policy, but may be a drawback in our setting, as the
 57 agent might not be consistently guided towards task completion. To address these shortcomings, we
 58 propose a method in which different states in the prior data are weighted different amounts by their
 59 estimated Q -value. More concretely, we propose a Q -weighted AIL approach that incentivizes the
 60 agent to move towards states in the prior data with higher value than its current state, so that agent
 61 may be guided consistently towards states closer to task completion.

62 Our contributions are as follows. First, we formalize the SLRL problem setting, which we believe to
 63 be a useful framework for modeling many situations in the real world. We next provide an intuitive
 64 argument and empirical analysis suggesting that learned reward shaping via distribution matching is
 65 better suited for this setting than finetuning without additional reward shaping. We identify challenges
 66 that uniquely arise in the SLRL setting with existing distribution matching approaches and propose
 67 a new approach, Q -weighted adversarial learning (QWALE), which is less sensitive to the quality
 68 of prior data available and provides the agent with a shaped reward towards completing the desired
 69 task a single time. Through our experiments, we explore the performance of different approaches in
 70 SLRL. We find that QWALE can meaningfully guide the agent to explore the state distribution in its
 71 prior data to complete the desired novel task 20-60% more successfully on four separate domains
 72 compared to existing distribution matching approaches and RL fine-tuning.

73 2 Related Work

74 **Autonomous RL.** In the context of deep RL, agents typically (but not always) are trained in episodic
 75 setting and are evaluated on the quality of the learned policy. Several recent works have developed
 76 algorithms that can learn without episodic resets [18, 6, 63, 42, 44, 15, 16, 24, 43]. Like our work,
 77 such methods aim to make it possible to learn without any episodic resets, but are typically still
 78 focused on acquiring an effective policy that can perform the task repeatedly, typically by training
 79 some auxiliary controller to enable the policy to “retry” the task multiple times without resets. In
 80 contrast, our aim is to develop an algorithm that can solve the task once, but as quickly as possible,
 81 which introduces a unique set of challenges as we discussed above.

82 **Continual RL.** There is a rich literature on reinforcement learning in the continuing setting [29,
 83 41, 47, 60, 37, 28, 55] that considers maximizing the average reward accumulated over an infinite
 84 horizon without episodic resets. Such works often also consider regret minimization as the objective.
 85 SLRL is closely related, and can be viewed as a special case where the agent has access to a prior
 86 offline dataset and aims to solve a single task as quickly as possible in a new domain. While the focus
 87 in continual learning is on general “lifelong” methods or on exploration, our focus is on effectively
 88 leveraging prior data in a setting that is meant to be reflective of real-world tasks (for example, in
 89 robotics).

90 **Leveraging offline data in online RL.** Learning expert policies given prior interaction data has been
 91 extensively studied in imitation learning [1, 38, 13], inverse RL [33, 8, 65, 66], RL for sparse reward
 92 settings [3, 32, 36, 20, 50] and offline RL [27, 26, 23, 56, 32]. Across all these diverse topics, the goal
 93 is to learn a competent policy that can solve the task efficiently whereas the objective in this work is
 94 to complete the task in a single trial as quickly as possible. To this end, we build on recent adversarial
 95 approaches to inverse RL [21, 11, 45, 49, 25, 64] to encourage agent’s state visitation towards expert
 96 prior data where the agent is likely to be successful. Prior methods have also studied adversarial
 97 inverse RL and imitation learning with non-expert data [51, 46, 57, 52, 54, 53, 5, 2]. However, as
 98 we discuss in Section 6, these approaches need to be adapted for the SLRL setting to be efficient at
 99 completing the task and to handle novelty, for example when the dynamics may have changed.

100 **Transfer and adaptation in RL.** Many prior works have studied the problem of adapting in presence
 101 of shifts between train and test settings, often in a specific problem setting such as sim2real trans-
 102 fer [40, 48, 35, 30] or fast adaptation via meta-learning [9, 34, 31, 67, 10]. A common theme in these
 103 works is that the algorithm can often train in preparation for adaptation at test-time, thus affecting the
 104 prior experiences it may collect. In contrast, the SLRL setting lays algorithmic emphasis on online
 105 exploration and adaptation, as the agent has access to fixed prior dataset of experiences. Other transfer
 106 learning approaches adapt the weights of the policy to a new environment or task, either through rapid
 107 zero-shot adaptation [19, 61] or through extended episodic online training [22, 39, 7, 59, 58]. Unlike
 108 the latter, we focus on adaptation within a single episode, but, unlike the former, with a focus on
 109 extended exploration and learning over tens of thousands of timesteps. This problem setting leads to
 110 unique challenges, namely that the agent must autonomously recover from mistakes, hence requiring
 111 a distinct approach.

112 3 Preliminaries

113 In this section, we describe some preliminaries before formalizing our problem statement in the
 114 following section. We consider an agent that operates in a Markov decision process (MDP) consisting
 115 of the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, \mathcal{R}, \rho, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the agent’s action space,
 116 $p(s_{t+1}|s_t, a_t)$ represents the environment’s transition dynamics, $\mathcal{R} : \mathcal{S} \rightarrow \mathbb{R}$ indicates the reward
 117 function, $\rho : \mathcal{S} \rightarrow \mathbb{R}$ denotes the initial state distribution, and $\gamma \in [0, 1)$ denotes the discount
 118 factor. In typical reinforcement learning, the objective is find a policy π that maximizes $J(\pi) =$
 119 $\mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t)]$.

120 Although our method is a reward-driven RL algorithm, we utilize concepts from imitation learning
 121 to overcome sparse rewards, utilizing potentially suboptimal prior data. To this end, we build on
 122 adversarial imitation learning (AIL), which uses prior data $\mathcal{D}_{\text{prior}}$ in the form of expert demonstrations
 123 (we will relax this requirement) to recover the expert’s policy. One such method is GAIL [21], which
 124 finds a policy π_{θ} that minimizes the Jensen-Shannon divergence between its stationary distribution
 125 and the expert data. It does so by training a discriminator network $D : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1)$, alternating
 126 updates with updates to the policy π . Concretely, D and π are learned by optimizing the following:
 127 $\min_{\pi} \max_{D \in (0,1)^{\mathcal{S} \times \mathcal{A}}} \mathbb{E}_{\pi} [\log(D(s, a))] + \mathbb{E}_{\pi_E} [\log(1 - D(s, a))] - \lambda H(\pi)$. In Section 5, we will see
 128 that AIL-style discriminator-based approaches can be adapted to the SLRL problem setting *without*
 129 demonstration data, and will therefore form the basis of our method.

130 4 Single-Life Reinforcement Learning

131 In this section, we formalize our problem setting, the *single-life reinforcement learning (SLRL)*
 132 problem. The defining characteristic of SLRL is that the agent is given a single “life”, i.e. trial, to
 133 complete a desired task, with the trial ending when the task is completed. The agent must complete
 134 the task autonomously, without access to any human interventions or resets.

135 In the real world, when faced with situations where a task must be completed once, an agent typically
 136 has some prior knowledge. E.g., an agent tasked with finding water on Mars may have experience

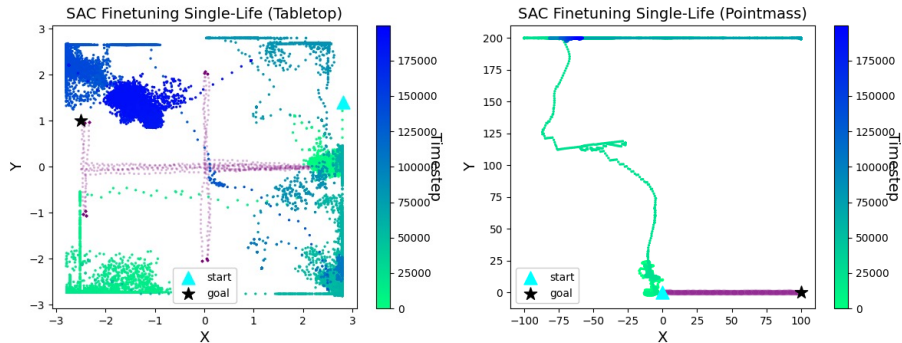


Figure 2: We visualize the online state visitation plots in the Tabletop and Pointmass environments of a single-life trial using SAC finetuning. We plot the location of the mug throughout the agent’s single life for the tabletop and the location of the agent for the pointmass, colored green to blue by timestep, along with expert demo states (purple). In both environments, the agent fails to recover from novel states and complete the task.

137 looking for water in a desert on Earth. We will therefore assume access to offline prior data of some
 138 sort that the agent may use for pretraining and during its single life. In many cases such as the Mars
 139 example, we may not have expert prior data of the desired task in the desired environment. Hence, we
 140 emulate this setting by providing the agent with prior data from a related environment and deploying
 141 its single life on a domain with a distribution shift.

142 We can formalize this setting as follows. We are given prior data $\mathcal{D}_{\text{prior}}$, which consists of transitions
 143 from some source MDP $\mathcal{M}_{\text{source}}$. The agent will then interact with a target MDP defined by $\mathcal{M}_{\text{target}} =$
 144 $(\mathcal{S}, \mathcal{A}, p, \mathcal{R}, \rho, \gamma)$. We assume that the target MDP has an aspect of novelty not present in the
 145 source MDP, such as different dynamics $p(s_{t+1} | s_t, a_t)$ or a different initial state distribution ρ .
 146 Naturally, the more similar the domains are, the easier the problem becomes, and the effectiveness of
 147 any algorithm will be strongly dependent on the degree of similarity, though formalizing a precise
 148 assumption on similarity between the source and target domain is difficult. The reward between the
 149 two MDPs is the same, meaning the agent is still trying to accomplish the same task in the target
 150 domain as in the source. The problem setting may be extended to include multiple source MDPs or a
 151 series of target MDPs.

152 The goal in SLRL is to accumulate as much reward as possible in a single trial in $\mathcal{M}_{\text{target}}$. Maintaining
 153 the same notation as the previous section, the SLRL problem aims to maximize $J = \sum_{t=0}^h \gamma^t \mathcal{R}(s_t)$,
 154 where h is the trial horizon, which may be ∞ . In general, we expect the task reward to be such that
 155 learning *only* from task rewards during the single life is difficult, for example because the reward is
 156 very sparse or even awarded only upon successful completion of the task (which can only happen
 157 once in the entire single life deployment). We assume that there are no sink states beyond a terminal
 158 success state, such that it is possible for the agent to autonomously recover from any mistake.

159 Note that this setup is essentially the same as the widely studied regret minimization problem in
 160 exploration [29]. However, while regret minimization is typically studied in the context of RL
 161 exploration theory, our aim with SLRL is to study a particular *special case* of the more general regret
 162 minimization framework that is meant to reflect a realistic setting in real-world RL (e.g., robotics)
 163 where an agent with prior experience must solve a task in a single (potentially long) trial.

164 As we analyze in Section 5, algorithms designed for episodic policy learning do not perform well
 165 in the SLRL problem setting, even when the policy and replay buffer are pre-trained and seeded
 166 with the prior data, because they do not quickly recover from mistakes to get back onto a good state
 167 distribution. In Section 6, we will discuss an approach based on distribution matching that attempts
 168 to address this issue.

169 5 Single-Life Performance of Online RL

170 Now that we have described the problem setting, we now empirically analyze online RL, which is
 171 designed for episodic learning, in the SLRL setting. Namely, we consider finetuning SAC [17], which
 172 pre-trains a policy and value function in the source setting and fine-tunes during the single trial in the
 173 target environment.

174 As we will see in Section 7, finetuning SAC performs poorly in the SLRL setting. We save the
 175 details of the experimental setup until Section 7, but to first motivate the use of distribution-matching
 176 approaches in our problem setting, we analyze the state visitation of SAC finetuning in the online
 177 phase for the Tabletop and Pointmass domains (see Figure 2). A key challenge of SLRL (and also
 178 fully autonomous RL in general) is that if the agent falls off of the distribution, it cannot rely on
 179 resets to get back on track. Since there is a gap between the source and target domains, the agent
 180 will inevitably find itself in states that are out of distribution from the prior data. A value function
 181 pre-trained on the source data could in principle be used to evaluate states and guide the agent back
 182 towards good states, but it will be inaccurate on states outside of the prior data [12]; then, when the
 183 value of some of those out-of-distributions states is overestimated, the value function may misguide
 184 the policy away from good states. Hence, fine-tuning a pre-trained value function via online RL will
 185 not explicitly encourage the agent to get back on distribution, especially in sparse reward envs. As a
 186 result, the agent may spend a lot of time (perhaps infinite time) drifting once it falls out of distribution,
 187 which we see occurs in Figure 2.

188 On the other hand, distribution-matching methods like GAIL [21] will explicitly encourage the
 189 agent to get back on distribution, by giving higher rewards on distribution than off distribution.
 190 However, existing adversarial imitation learning methods assume that the prior data consists of
 191 expert demonstrations, and they train the agent to match the entire demo distribution, which is not
 192 necessarily the ideal distribution to match. In the following section, we will discuss a method that
 193 aims to address these shortcomings.

194 6 Q -weighted Adversarial Learning (QWALE)

195 In this section, we will present our method for addressing SLRL, which we call QWALE. The key
 196 insight in QWALE is to utilize the prior data $\mathcal{D}_{\text{prior}}$ to handle the sparse and uninformative reward
 197 information in the target domain. Our first observation is that the framework of AIL already provides
 198 a reasonable starting point, though it is not sufficient by itself: rather than using only the task reward,
 199 which is too sparse to be useful in a single episode, we can bias the agent to seek out states that
 200 are similar to those seen in the prior data. However, since our goal is not to learn a policy that
 201 repeatedly performs the task, but rather to solve it as quickly as possible once, we do not actually
 202 want to *learn to imitate* prior data, but rather to seek out states that resemble the *best* states in the
 203 prior data, with better states being more preferred. This is especially important when the prior data is
 204 not actually optimal, but might consist of arbitrarily suboptimal states. We will discuss how this can
 205 be accomplished with a modification of AIL which, instead of treating prior data as equally desirable,
 206 preferentially drives the agent toward states that resemble the best states in the prior data.

207 6.1 Algorithm Description

208 In SLRL, we may have shifts in dynamics online, in which case matching state-action distributions as
 209 done in GAIL ([21]) may not be appropriate. GAIL with a state discriminator will help lead the agent
 210 back to the prior data distribution but by matching the entire demo state distribution, the discriminator
 211 does not necessarily incentivize the agent to go towards task completion. Our algorithm’s desired
 212 behavior is to lead the agent to nearby states within distribution of the prior data if it is out of
 213 distribution and to nearby states closer to task completion if in distribution. Our proposed method
 214 for shaping relies on the intuition that as states gradually get closer to task completion, even within
 215 expert data, they should have gradually higher values as well.

216 We propose Q -weighted adversarial learning (QWALE), which trains a Q -weighted discriminator. In
 217 order to use the prior data effectively, we use a fixed Q -function $Q(s, a)$ trained in the source MDP
 218 to distinguish between useful transitions and ones that may be less useful. This Q -function may be
 219 obtained through RL pretraining, which is what we use for our experiments, or a variety of other
 220 ways, such as offline RL or Monte Carlo estimation. We train the discriminator in a similar manner
 221 as GAIL, where the positives come from the offline data and negatives from online experience. To
 222 take into account the varied quality of the data, we use the intuition that the closer a state is to task
 223 completion, the higher its value should be. In particular, a state in the prior data should get smaller
 224 weight if it has worse value than the agent’s current state, so the agent is consistently incentivized
 225 to move towards states in the prior data with higher value than its current state. Therefore, when
 226 training the discriminator, we weight the positive states s by $\exp(Q(s, a) - b)$ and the negatives by
 227 $\exp(-Q(s, a) + b)$, where b is an implementation detail, which we discuss in the Appendix. We
 228 normalize the Q -values to be between 0 and 1 and train the Q -weighted discriminator in alternating

229 updates with SAC updates that finetune the policy and critic in an AIL fashion. In this manner,
 230 QWALE extends AIL to the general setting with any prior data.

231 The goal of our weighted discriminator training procedure is to obtain a discriminator that, when
 232 used as a reward, will drive the agent toward states that it believes would lead to better outcomes
 233 than its present state, based on the prior data. The Q -function quantifies the agent’s belief from the
 234 prior data that a particular state will lead to high reward, making this a natural choice for estimating
 235 how desirable a state is at any given time. Hence, using Q -values to weight the examples for the
 236 discriminator will cause the discriminator to prefer states that are closer to the goal over states that
 237 are further away. This is significantly different from the behavior we would expect to see if we were
 238 simply imitating the optimal policy, as this would give equal weight to all of the transitions along an
 239 optimal path. When the reward is more complex, using Q -values as weights generalizes this intuition.

240 6.2 Practical Implementation

Algorithm 1 Q-WEIGHTED ADVERSARIAL LEARNING (QWALE)

```

1: // Single Trial Deployment
2: Require:  $\mathcal{D}_{\text{prior}}$ , test MDP  $\mathcal{M}_{\text{test}}$ , pretrained critic  $Q(s, a)$ , and (optionally) policy  $\pi$ ;
3: Initialize: replay buffer for online transitions  $\mathcal{D}_{\text{online}}$ ; parameters  $\phi$  for discriminator  $q_{\phi}(\text{prior} \mid s_t)$ , timestep
    $t = 0$ 
4: while task not complete do
5:   Sample  $a \sim \pi(\cdot \mid s_t)$ 
6:    $\phi \leftarrow \phi - \eta \nabla_{\phi} L(\phi)$  // Update discriminator according to Eq. 1
7:    $r'(s_t) = r(s_t) + \log q_{\phi}(\text{prior} \mid s_t)$ 
8:    $Q(s, a), \pi \leftarrow \text{SAC}(Q(s, a), \pi, \mathcal{D}_{\text{prior}} \cup \mathcal{D}_{\text{online}}, r')$ 
9:   Increment  $t$ 

```

241 We optimize our objective using maximum entropy off-policy RL with the SAC algorithm ([17]),
 242 modified in a similar manner as we did with GAIL in the previous section. In particular, we learn
 243 an additional discriminator $q_{\theta}(\text{prior} \mid s_t)$, optimized using standard cross-entropy loss, which is
 244 weighted accordingly:

$$L(\phi) = -\mathbb{E}_{\mathcal{D}_{\text{prior}}} [\exp(Q(s, a) - b) \log q_{\phi}(\text{prior} \mid s)] - \mathbb{E}_{\mathcal{D}_{\text{online}}} [\exp(-Q(s, a) + b) \log q_{\phi}(\text{online} \mid s)]. \quad (1)$$

245 The discriminator is used to modify the rewards when updating off-policy from all experience—prior
 246 and online. At single trial test time, the actor and critic are optionally initialized with the pretrained
 247 weights, and the replay buffer is initialized with the offline data. For details such as network
 248 architecture and hyperparameters, see Appendix A. We present the full algorithm in Algorithm 1.

249 7 Experiments

250 The goal of our experiments is to answer the following questions: (1) How does QWALE compare to
 251 prior reinforcement learning and distribution matching approaches in single-life RL settings? (2) Do
 252 distribution matching approaches help agents learn to recover from novel situations in single-life RL?
 253 (3) How does QWALE compare to different variants of adversarial imitation learning, with different
 254 prior datasets?

255 7.1 Experimental Setup

256 To answer the above questions, we construct four single-life RL domains with varying prior datasets
 257 and sources of novelty, and then measure performance both in terms of speed of task completion and
 258 overall single-life success. In this subsection, we describe this experimental set-up in detail.

259 **Environments.** We consider the following four problem domains. First, in the Tabletop-Organization
 260 environment from the EARL benchmark [43], the agent is tasked with bringing a mug to one of four
 261 different locations designated by a goal coaster. The prior data always has the same starting position
 262 of the mug. In the target environment, the starting position is in a new location unseen in the prior
 263 data. Second, the Pointmass setting tasks an agent to move in 2D from its starting location at the
 264 origin $(0, 0)$ to the point $(100, 0)$. The target environment introduces a dynamics shift in the form of a
 265 strong “wind”, where the agent is involuntarily pushed upward in the y-direction each step. Third, we
 266 construct a modified HalfCheetah environment, in which it is difficult but feasible for the cheetah to
 267 recover when flipped over. The target environment includes hurdles that the cheetah must jump over,

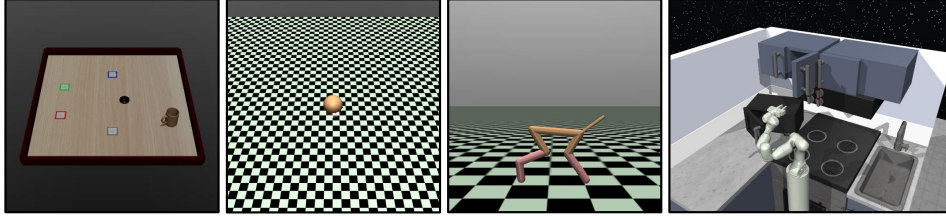


Figure 3: We evaluate in four different domains, including Tabletop-Organization, Pointmass, HalfCheetah, and a Franka-Kitchen environment with a microwave and cabinet. At test time, an aspect of novelty is introduced in each environment—new initial mug positions for Tabletop, wind for Pointmass, hurdles for the HalfCheetah, and a new combination of tasks for the Franka-Kitchen.

268 as the prior data does not include these obstacles. Finally, we evaluate on a modified Franka-Kitchen
 269 environment, adapted from [14], where the task is to close a microwave and a hinged cabinet. The
 270 prior data only contains trajectories of closing the microwave and the hinged cabinet separately, so
 271 the agent must figure out online how to complete both tasks in a row. In other words, both objects are
 272 open at the start of single-life RL, and the agent has only previously seen instances where only one is
 273 open. For the latter two environments, dense rewards are given, and the discriminator-based reward is
 274 added to the extrinsic reward during single-life training. These environments are shown in Figure 3.
 275 Further details on the environments are given in Appendix A.

276 **Comparisons.** To answer question (1), we compare QWALE to three alternative methods: (a) SAC
 277 fine-tuning, which pre-trains a policy and value function in the source setting and fine-tunes for
 278 a single, long episode in the target environment, (b) SAC-RND, which additionally includes an
 279 RND exploration bonus [4] during single-life fine-tuning, and (c) GAIL-s, which runs generative
 280 adversarial imitation learning [21, 25] where the discriminator only operates on the current state
 281 s . We choose for the discriminator to only look at s so that it is less susceptible to dynamics shift
 282 between the source data and target environment. We additionally compare to GAIL-sa, which passes
 283 both the current state and action to the discriminator. All methods use soft actor-critic (SAC) [17] as
 284 the base RL algorithm.

285 **Prior datasets.** For all four environments, we evaluate SLRL using data collected through RL as our
 286 prior data. More specifically, we run SAC in the source MDP in the standard episodic RL setting for
 287 K steps and take the last 50,000 transitions as the prior data. K is chosen such that the prior data
 288 contains some good transitions but has not converged to an optimal policy yet. While we are able to
 289 run episodic RL in the source MDP, this is not a requirement for SLRL, as long as prior data in the
 290 source MDP is available. For all methods, including QWALE, GAIL variants, and SAC fine-tuning
 291 variants, the policy and value function are pretrained in this manner for the initialization of single-life
 292 RL. We note that AIL methods like GAIL typically assume that the prior data consists of expert
 293 demonstrations but we apply the algorithm only using mixed quality prior data, unless otherwise
 294 noted. In particular, Section 7.4 further evaluates AIL methods using demos as prior data, using 10
 295 demonstrations for the Tabletop environment and 3 demonstrations for the Pointmass domain. We
 296 include such experiments to answer question (3), i.e. to investigate how the quality of prior data may
 297 affect performance.

298 **Evaluation Metrics.** To evaluate each method in each environment, we report the average and median
 299 number of steps taken before task completion across 10 seeds along with the standard error and
 300 success rate (out of 10). During single-life RL, for all environments, the agent is given a maximum of
 301 200,000 steps to complete the task. If it has not completed the task after 200k steps, then 200k is
 302 logged as the total number of steps, and the run is marked as unsuccessful.

303 7.2 Results using mixed data as prior data

304 In this subsection, we aim to answer our first experimental question and study how QWALE performs
 305 compared to prior reinforcement learning and distribution matching approaches in single-life RL
 306 settings. As seen in Figure 4 and Table 1, we find that QWALE achieves the lowest average and
 307 median number of steps as well as highest number of successes on three out of the four domains,
 308 and performs comparably to the other methods on the fourth environment, Franka-Kitchen. On the
 309 Tabletop and Pointmass environments, QWALE takes less than half of the number of steps on average
 310 as the next best performing method. It is possible that the method does not work as well on the
 311 Franka-Kitchen environment because the pretrained Q -function may be quite inaccurate, as there is

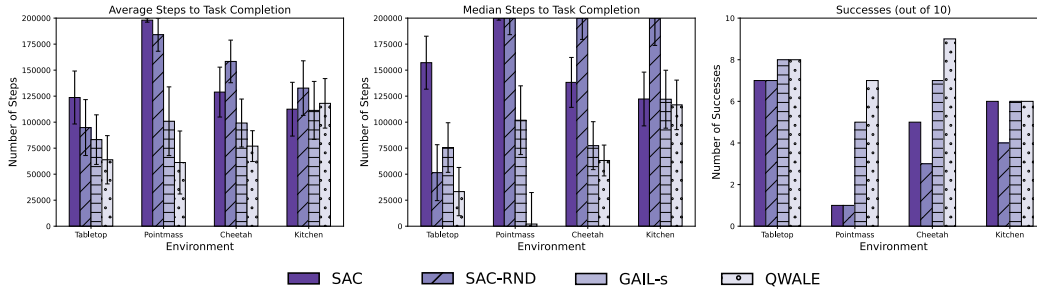


Figure 4: We evaluate the performance of QWALE to finetuning SAC and GAIL in our four environments using mixed data collected through RL as prior data. We omit the results of Behavior Cloning (BC) in the plots, as it is unsuccessful at completing the task in every domain due to the distribution shift. We plot the average and median number of steps to task completion along with the number of successes, taken over 10 seeds. We find that GAIL outperforms SAC in 3 out of 4 domains, and QWALE significantly outperforms GAIL on 3 out of 4 domains and performs comparably on the fourth.

	Method	Avg \pm Std error	Success / 10	Median		Method	Avg \pm Std error	Success / 10	Median
Tabletop	GAIL-s	83.2k \pm 23.8k	8	75.6k	Cheetah	GAIL-s	99.2k \pm 23.0k	7	77.4k
	GAIL-sa	61.5k \pm 28.7k	7	2.4k		GAIL-sa	102.0k \pm 19.3k	8	85.6k
	QWALE (ours)	23.1k \pm 7.9k	8	15.5k		QWALE (ours)	77.0k \pm 14.8k	9	63.2k
Pointmass	GAIL-s	100.9k \pm 33.0k	5	101.9k	Kitchen	GAIL-s	111.3k \pm 27.9k	6	122.1k
	GAIL-sa	140.4k \pm 30.3k	3	200.0k		GAIL-sa	127.8k \pm 26.9k	5	189.1k
	QWALE (ours)	61.2k \pm 30.2k	7	2.1k		QWALE (ours)	118.1k \pm 23.8k	6	116.6k

Table 1: Discriminator-based Approaches on Mixed Data. We see that in each of the four experimental domains, the three QWALE methods typically outperform GAIL and GAIL-sa on the average number of steps needed before task completion, and on 3 out of the 4 environments, QWALE substantially improves performance over both GAIL variants. All methods are evaluated over 10 runs.

312 a global distribution shift in the state space at test time—the agent has never seen both objects open
 313 before in the prior data. GAIL also outperforms finetuning SAC across three of the four domains.
 314 These results demonstrate the suitability of distribution-matching approaches over RL finetuning in
 315 the SLRL setting. We see that guidance particularly towards a good state distribution is important,
 316 as we compare to finetuning SAC with an exploration bonus through random network distillation
 317 ([4]). From Figure 4, although RND may improve performance, particularly in the Tabletop domain,
 318 it generally does not perform as well as the distribution matching approaches, especially QWALE,
 319 showing that simply increasing exploration is not enough. Furthermore, these results show that the
 320 additional shaping provided by weighting the prior data by Q -value when training the discriminator
 321 can significantly improve guidance towards the goal. While GAIL gives equal weight to all transitions
 322 along an optimal path, the agent in QWALE is consistently guided towards states in the prior data
 323 with higher Q -value, leading to more efficient and reliable single-life task completion.

324 7.3 Analysis of distribution matching approaches

325 Next, to answer question (2), we analyze how QWALE helps agents learn to recover from novel
 326 situations in SLRL. To do so, we visualize QWALE’s state visitation in the Tabletop and Pointmass
 327 domains throughout a single lifetime. We color the trajectories according to timestep for both
 328 methods as well as by reward (discriminator score). The coloring in the timestep-colored plots is
 329 highly correlated with that in the reward-colored plots, showing how the reward gradually guides the
 330 agent towards the goal. In particular, when the agent is out of distribution, the agent is incentivized to
 331 explore states that will lead it closer back to the prior state distribution, and when the agent is within
 332 distribution, it is incentivized to move to states closer to the goal, leading to efficient task completion.

333 7.4 Using demos as prior data

334 Finally, we evaluate the performance of different discriminator-based approaches in the two SLRL
 335 problem settings—Tabletop and Pointmass—where demonstration data is available as prior data. We
 336 compare using AIL with a state-only discriminator (GAIL-s) as well as with a state-action discrimina-
 337 tor (GAIL-sa) to our proposed Q -weighted discriminator method (QWALE). With the latter method,
 338 we have access to the same Q -function pretrained when collecting prior data using standard RL for
 339 the mixed data experiments above, but we do not initialize any of the algorithms at test time with the
 340 pretrained policy and critic weights.

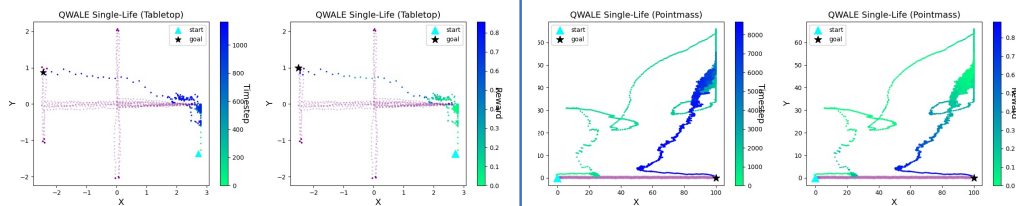


Figure 5: We visualize the online state visitation plots in the Tabletop (left) and Pointmass (right) environments of a single-life trial using QWALE. We plot the location of the mug throughout the agent’s single life for the tabletop and the location of the agent for the pointmass as well as expert demo states (purple). We color the trajectories green to blue according to timestep as well by reward (discriminator score) for the distribution-matching approach.

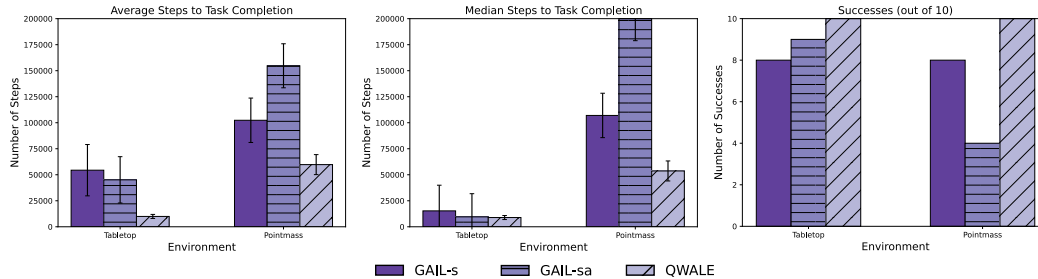


Figure 6: Discriminator-based Approaches using Expert Demo Data. Given demonstration data, in both the Tabletop and Pointmass domains, QWALE significantly outperforms both GAIL variants in almost all metrics.

341 From Figure 6, the GAIL variants are both able to consistently solve the task in the Tabletop
 342 domain, but unsurprisingly, GAIL-sa does especially poorly in the Pointmass domain, where the
 343 dynamics have changed at test time. Compared to the two GAIL variants, QWALE gives a significant
 344 improvement in both domains. These results demonstrate how more detailed reward shaping towards
 345 the completion of the desired task can be helpful in the SLRL setting even with demonstrations as
 346 prior data. Moreover, comparing these results with those in Table 1, while access to expert data
 347 as prior data unsurprisingly improves the performance of GAIL methods, it can also improve the
 348 performance of QWALE.

349 8 Conclusion

350 In this paper, we formalized and studied a problem setting underlying single-life reinforcement
 351 learning: settings where an agent needs to autonomously complete a task once while drawing upon
 352 prior experience from a related environment. We found that standard fine-tuning via RL is ill-suited
 353 for this problem because the algorithm struggles to recover from mistakes and novel situations. We
 354 hypothesized that this observation stems from the fact that resets in episodic RL prevent algorithms
 355 from needing to recover, whereas single-life RL and continuing settings in general do demand the
 356 agent to find its way back to good states on its own. We then postulated that distribution matching
 357 methods that aim to match the distribution of related prior data may help agents recover via reward
 358 shaping, and presented a new distribution matching method, QWALE, that weights examples by their
 359 Q -value. Our experiments verified that distribution matching approaches indeed do make better use
 360 of prior data, and that QWALE is competitive with or outperforms prior distribution matching methods
 361 on four single-life RL problems.

362 While QWALE can efficiently complete novel target tasks in a single episode without any interventions,
 363 important limitations remain. No algorithm, including QWALE, was able to complete the target task
 364 with 100% success, indicating that future works should aim to improve an algorithm’s ability to
 365 solve tasks consistently. Moreover, the methods that we evaluated all used a pre-trained policy and
 366 value function from the source domain, which may be difficult to obtain in some source scenarios,
 367 as opposed to only obtaining some demonstrations or offline data. Finally, it would be interesting
 368 to explore problems with greater degrees of novelty between the source and target environments.
 369 We expect that such settings would place even greater importance on autonomy and exploration,
 370 requiring sophisticated strategies for both recovering to known states and exploring new strategies.
 371 By publicly releasing and open-sourcing the environments and code upon publication, we hope that
 372 future work can more easily explore these interesting questions and continue to make progress on
 373 allowing RL agents to autonomously complete tasks within a single lifetime.

374 References

- 375 [1] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot
376 learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- 377 [2] Mark Beliaev, Andy Shih, Stefano Ermon, Dorsa Sadigh, and Ramtin Pedarsani. Imitation
378 learning by estimating expertise of demonstrators. *arXiv preprint arXiv:2202.01288*, 2022.
- 379 [3] Tim Brys, Anna Harutyunyan, Halit Bener Suay, Sonia Chernova, Matthew E Taylor, and
380 Ann Nowé. Reinforcement learning from demonstration through shaping. In *Twenty-fourth
381 international joint conference on artificial intelligence*, 2015.
- 382 [4] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random
383 network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- 384 [5] Zhangjie Cao, Zihan Wang, and Dorsa Sadigh. Learning from imperfect demonstrations via
385 adversarial confidence transfer. *arXiv preprint arXiv:2202.02967*, 2022.
- 386 [6] Benjamin Eysenbach, Shixiang Gu, Julian Ibarz, and Sergey Levine. Leave no trace: Learning
387 to reset for safe and autonomous reinforcement learning. *arXiv preprint arXiv:1711.06782*,
388 2017.
- 389 [7] Benjamin Eysenbach, Swapnil Asawa, Shreyas Chaudhari, Sergey Levine, and Ruslan Salakhut-
390 dinov. Off-dynamics reinforcement learning: Training for transfer with domain classifiers.
391 *arXiv preprint arXiv:2006.13916*, 2020.
- 392 [8] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal
393 control via policy optimization. In *International conference on machine learning*, pages 49–58.
394 PMLR, 2016.
- 395 [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-
396 tation of deep networks. In *International conference on machine learning*, pages 1126–1135.
397 PMLR, 2017.
- 398 [10] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning.
399 In *International Conference on Machine Learning*, pages 1920–1930. PMLR, 2019.
- 400 [11] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse
401 reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 402 [12] Justin Fu, Aviral Kumar, Matthew Soh, and Sergey Levine. Diagnosing bottlenecks in deep
403 q-learning algorithms. In *International Conference on Machine Learning*, pages 2021–2030.
404 PMLR, 2019.
- 405 [13] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence mini-
406 mization perspective on imitation learning methods. In *Conference on Robot Learning*, pages
407 1259–1277. PMLR, 2020.
- 408 [14] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay
409 policy learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv
410 preprint arXiv:1910.11956*, 2019.
- 411 [15] Abhishek Gupta, Justin Yu, Tony Z Zhao, Vikash Kumar, Aaron Rovinsky, Kelvin Xu, Thomas
412 Devlin, and Sergey Levine. Reset-free reinforcement learning via multi-task learning: Learning
413 dexterous manipulation behaviors without human intervention. In *2021 IEEE International
414 Conference on Robotics and Automation (ICRA)*, pages 6664–6671. IEEE, 2021.
- 415 [16] Abhishek Gupta, Corey Lynch, Brandon Kinman, Garrett Peake, Sergey Levine, and Karol
416 Hausman. Bootstrapped autonomous practicing via multi-task reinforcement learning. *arXiv
417 preprint arXiv:2203.15755*, 2022.
- 418 [17] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
419 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International
420 conference on machine learning*, pages 1861–1870. PMLR, 2018.
- 421 [18] Weiqiao Han, Sergey Levine, and Pieter Abbeel. Learning compound multi-step controllers
422 under unknown dynamics. In *2015 IEEE/RSJ International Conference on Intelligent Robots
423 and Systems (IROS)*, pages 6435–6442. IEEE, 2015.
- 424 [19] Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà, Pieter Abbeel, Alexei A Efros, Lerrel
425 Pinto, and Xiaolong Wang. Self-supervised policy adaptation during deployment. *arXiv preprint
426 arXiv:2007.04309*, 2020.

- 427 [20] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
428 John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations. In
429 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 430 [21] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
431 *information processing systems*, 29, 2016.
- 432 [22] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual rein-
433 forcement learning: A review and perspectives. *arXiv preprint arXiv:2012.13490*, 2020.
- 434 [23] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:
435 Model-based offline reinforcement learning. *Advances in neural information processing systems*,
436 33:21810–21823, 2020.
- 437 [24] Jigang Kim, J hyeon Park, Daesol Cho, and H Jin Kim. Automating reinforcement learning
438 with example-based resets. *IEEE Robotics and Automation Letters*, 2022.
- 439 [25] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan
440 Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in
441 adversarial imitation learning. *arXiv preprint arXiv:1809.02925*, 2018.
- 442 [26] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning
443 for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:
444 1179–1191, 2020.
- 445 [27] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:
446 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 447 [28] Vincenzo Lomonaco, Karan Desai, Eugenio Culurciello, and Davide Maltoni. Continual
448 reinforcement learning in 3d non-stationary environments. In *Proceedings of the IEEE/CVF*
449 *Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- 450 [29] Sridhar Mahadevan. Average reward reinforcement learning: Foundations, algorithms, and
451 empirical results. *Machine learning*, 22(1):159–195, 1996.
- 452 [30] Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active
453 domain randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- 454 [31] Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine,
455 and Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-
456 reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.
- 457 [32] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Over-
458 coming exploration in reinforcement learning with demonstrations. In *2018 IEEE international*
459 *conference on robotics and automation (ICRA)*, pages 6292–6299. IEEE, 2018.
- 460 [33] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*,
461 volume 1, page 2, 2000.
- 462 [34] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms.
463 *arXiv preprint arXiv:1803.02999*, 2018.
- 464 [35] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real
465 transfer of robotic control with dynamics randomization. In *2018 IEEE international conference*
466 *on robotics and automation (ICRA)*, pages 3803–3810. IEEE, 2018.
- 467 [36] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
468 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
469 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 470 [37] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experi-
471 ence replay for continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc,
472 E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32.
473 Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/
474 file/fa7cdfad1a5aaf8370ebeda47a1ff1c3-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/fa7cdfad1a5aaf8370ebeda47a1ff1c3-Paper.pdf).
- 475 [38] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and
476 structured prediction to no-regret online learning. In *Proceedings of the fourteenth interna-*
477 *tional conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and
478 Conference Proceedings, 2011.

- 479 [39] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick,
480 Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv*
481 *preprint arXiv:1606.04671*, 2016.
- 482 [40] Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real single-image flight without a single real
483 image. *arXiv preprint arXiv:1611.04201*, 2016.
- 484 [41] Anton Schwartz. A reinforcement learning method for maximizing undiscounted rewards. In
485 *Proceedings of the tenth international conference on machine learning*, volume 298, pages
486 298–305, 1993.
- 487 [42] Archit Sharma, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Au-
488 tonomous reinforcement learning via subgoal curricula. In M. Ranzato, A. Beygelz-
489 imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neu-
490 ral Information Processing Systems*, volume 34, pages 18474–18486. Curran Asso-
491 ciates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/file/
492 99c83c904d0d64fbef50d919a5c66a80-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/99c83c904d0d64fbef50d919a5c66a80-Paper.pdf).
- 493 [43] Archit Sharma, Kelvin Xu, Nikhil Sardana, Abhishek Gupta, Karol Hausman, Sergey Levine,
494 and Chelsea Finn. Autonomous reinforcement learning: Formalism and benchmarking. *arXiv*
495 *preprint arXiv:2112.09605*, 2021.
- 496 [44] Archit Sharma, Rehaan Ahmad, and Chelsea Finn. A state-distribution matching approach to
497 non-episodic reinforcement learning. *arXiv preprint arXiv:2205.05212*, 2022.
- 498 [45] Avi Singh, Larry Yang, Kristian Hartikainen, Chelsea Finn, and Sergey Levine. End-to-end
499 robotic reinforcement learning without reward engineering. *arXiv preprint arXiv:1904.07854*,
500 2019.
- 501 [46] Mingfei Sun and Xiaojuan Ma. Adversarial imitation learning from incomplete demonstrations.
502 *arXiv preprint arXiv:1905.12310*, 2019.
- 503 [47] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,
504 2018.
- 505 [48] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel.
506 Domain randomization for transferring deep neural networks from simulation to the real world.
507 In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages
508 23–30. IEEE, 2017.
- 509 [49] Faraz Torabi, Garrett Warnell, and Peter Stone. Adversarial imitation learning from state-only
510 demonstrations. In *Proceedings of the 18th International Conference on Autonomous Agents*
511 *and MultiAgent Systems*, pages 2229–2231, 2019.
- 512 [50] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas
513 Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations
514 for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint*
515 *arXiv:1707.08817*, 2017.
- 516 [51] Qing Wang, Jiechao Xiong, Lei Han, Han Liu, Tong Zhang, et al. Exponentially weighted
517 imitation learning for batched historical data. *Advances in Neural Information Processing*
518 *Systems*, 31, 2018.
- 519 [52] Ruohan Wang, Carlo Ciliberto, Pierluigi Amadori, and Yiannis Demiris. Support-weighted
520 adversarial imitation learning. *arXiv preprint arXiv:2002.08803*, 2020.
- 521 [53] Yunke Wang, Chang Xu, and Bo Du. Robust adversarial imitation learning via adaptively-
522 selected demonstrations. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International*
523 *Joint Conference on Artificial Intelligence, IJCAI-21*, pages 3155–3161. International Joint
524 Conferences on Artificial Intelligence Organization, 2021.
- 525 [54] Yunke Wang, Chang Xu, Bo Du, and Honglak Lee. Learning to weight imperfect demonstrations.
526 In *International Conference on Machine Learning*, pages 10961–10970. PMLR, 2021.
- 527 [55] Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, Hiteshi Sharma, and Rahul Jain. Model-
528 free reinforcement learning in infinite-horizon average-reward Markov decision processes. In
529 Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference*
530 *on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages
531 10170–10180. PMLR, 13–18 Jul 2020.

- 532 [56] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement
533 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 534 [57] Yueh-Hua Wu, Nontawat Charoenphakdee, Han Bao, Voot Tangkaratt, and Masashi Sugiyama.
535 Imitation learning from imperfect demonstration. In *International Conference on Machine*
536 *Learning*, pages 6818–6827. PMLR, 2019.
- 537 [58] Annie Xie and Chelsea Finn. Lifelong robotic reinforcement learning by retaining experiences.
538 *arXiv preprint arXiv:2109.09180*, 2021.
- 539 [59] Annie Xie, James Harrison, and Chelsea Finn. Deep reinforcement learning amidst lifelong
540 non-stationarity. *arXiv preprint arXiv:2006.10701*, 2020.
- 541 [60] Ju Xu and Zhanxing Zhu. Reinforced continual learning. In S. Bengio,
542 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, edi-
543 tors, *Advances in Neural Information Processing Systems*, volume 31. Curran As-
544 sociates, Inc., 2018. URL [https://proceedings.neurips.cc/paper/2018/file/](https://proceedings.neurips.cc/paper/2018/file/cee631121c2ec9232f3a2f028ad5c89b-Paper.pdf)
545 [cee631121c2ec9232f3a2f028ad5c89b-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/cee631121c2ec9232f3a2f028ad5c89b-Paper.pdf).
- 546 [61] Takuma Yoneda, Ge Yang, Matthew R Walter, and Bradly Stadie. Invariance through inference.
547 *arXiv preprint arXiv:2112.08526*, 2021.
- 548 [62] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond
549 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 550 [63] Henry Zhu, Justin Yu, Abhishek Gupta, Dhruv Shah, Kristian Hartikainen, Avi Singh, Vikash
551 Kumar, and Sergey Levine. The ingredients of real-world robotic reinforcement learning. *arXiv*
552 *preprint arXiv:2004.12570*, 2020.
- 553 [64] Zhuangdi Zhu, Kaixiang Lin, Bo Dai, and Jiayu Zhou. Off-policy imitation learning from
554 observations. In *the Thirty-fourth Annual Conference on Neural Information Processing Systems*
555 *(NeurIPS 2020)*, 2020.
- 556 [65] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy
557 inverse reinforcement learning. In *AAAI*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- 558 [66] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of
559 maximum causal entropy. In *ICML*, 2010.
- 560 [67] Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann,
561 and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-
562 learning. *arXiv preprint arXiv:1910.08348*, 2019.

563 A Appendix

564 A.1 Implementation Details and Hyperparameters

565 In our experiments, we use soft actor-critic [17] as our base RL algorithm. We use default hyperpa-
566 rameter values: a learning rate of $3e-4$ for all networks, optimized using Adam, with a batch size
567 of 256 sampled from the entire replay buffer (both prior and online data), a discount factor of 0.99.
568 The policy and critic networks are MLPs with 2 fully-connected hidden layers of size 256. For all
569 methods training a discriminator, it is parameterized as an MLP with 1 fully-connected hidden layer
570 of size 128 and trained with a batch size of 512. During the online trial, 1000 steps are taken as
571 initial collection steps before network updates begin. For all methods training a discriminator, we use
572 mixup regularization [62] to reduce the brittleness of the discriminator.

573 Following [43], we use a biased TD update, where $Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma Q(s_{t+1}, a_{t+1})$ if t is
574 not a multiple of 100, and $Q(s_t, a_t) \leftarrow r(s_t, a_t)$ if it is. We use this update for all our evaluated
575 methods online in order to improve stability. Since the online trial of single-life RL may have a large
576 training horizon with hundreds of thousands of steps, this may lead to unstable bootstrapping, as for
577 each t , $Q(s_t, a_t)$ bootstraps on $Q(s_{t+1}, a_{t+1})$. Following [44], for the auxiliary reward given by the
578 discriminator D , we use $r(s, a) = -\log(1 - D(s))$ instead of $r(s, a) = \log D(s)$ to further improve
579 stability.

580 For QWALE, the weighting of states is offset by a value b . This value may be treated like a constant
581 hyperparameter and tuned. Adding this value changes the bias on the discriminator, which in effect
582 adds a constant to the reward, though that constant changes over the course of training. In practice,
583 to avoid having to tune b , we just use the value of the most recent state as b , i.e. $b = Q(s_t, a_t)$. To
584 better interpret this value, with this weighting, b is a baseline value capturing some notion of current
585 progress. Prior data tends to get small weights if they have worse value than the current state, so the
586 agent is consistently incentivized to move towards states with higher value than its current state.

587 For all experiments using prior data collected through RL, the agent was initialized at test time
588 with the pretrained policy and critic. For QWALE, a copy of that critic was frozen and used when
589 calculating the weights for discriminator training. For all of the experiments with demonstration data
590 in Section 7.4, the policy and critic were not initialized with any pretrained weights.

591 A.2 Environment & Evaluation Details

592 *Tabletop-Organization.* The details for this environment are in [43]. The state space consists of the
593 gripper’s (x, y) position, the mug’s (x, y) position, the gripper’s state (whether attached to the mug or
594 not), and the current goal, for a total of 12 dimensions. The action space is 3 dimensional, consisting
595 of a delta in the gripper’s (x, y) position as well as an automatic gripper that will attach to the mug if
596 the gripper is close enough. The tabletop extends from -2.8 to 2.8 in both the x and y directions. In
597 the prior data, which consists either of 10 demonstrations or 50000 transitions collected through RL
598 after 350000 steps of training, the initial state always places the mug at position (2.5, 0.0), and the
599 goal is to place the mug at one of the following locations: (-2.5, -1.0), (-2.5, 1.0), (0, 2.0), (0, -2).
600 For the online trial when evaluating SLRL, the mug is placed either at (2.7, 1.5) or (2.7, -1.5) with
601 additional uniform randomness between (-0.15, 0.15) in both directions. This environment is also
602 goal-conditioned at test time and the goal is randomly set to be either (-2.5, -1.0) or (-2.5, 1.0). The
603 reward is 1 when the mug is within 0.15 distance of its goal position (at which point the single life
604 ends) and 0 everywhere else.

605 *Pointmass.* The Pointmass environment has a 6-dimensional state space consisting of the agent’s
606 (x, y) position, its (x, y) velocity, and the (x, y) coordinates of the goal. The environment extends
607 between -100 and 100 along the x axis and between -200 and 200 along the y axis. The action space
608 is 2-D, consisting of the delta in both directions, clipped between -1 and 1 for a single action. The
609 prior data consists of 3 demonstrations or 50000 transitions collected through RL after 350000 steps
610 of training. The agent starts at (0, 0) and the goal is at (100, 0) for the prior data and online trial.
611 During the online trial, a strong “wind” is introduced, where a random amount between 0.8 and 0.9 is
612 added to the agent’s y coordinate and 0.2 is subtracted from the agent’s x coordinate at each step.
613 The reward is 1 when the agent is within a distance of 2 of the goal position and 0 everywhere else.

614 *HalfCheetah.* The HalfCheetah environment has a state space with 18 dimensions, consisting of the
615 position and velocity of each joint. The prior data consists of 50000 transitions collected through RL
616 after 150000 steps of training. The reward is $r_t = \Delta x_t - 0.1 * ||a_t||_2^2$. At test time, 10 hurdles are

617 included in the environment, spread between the x-coordinate of 7 and 260. The cheetah starts at 0
618 and its single life is considered successful when it gets to the coordinate 300, although the information
619 about the hurdles or goal are not included in the state space.

620 *Franka-Kitchen*. The Franka-Kitchen is adapted from [14, 43]. The state space consists of a 9 DoF
621 position-controlled Franka-robot with a microwave and hinged cabinet. The prior data consists of
622 50000 transitions collected through standard episodic RL after 950000 steps of training, where one of
623 the microwave or cabinet is open, and the task is to close that object. At test time, both are open, and
624 the task is to close both objects. The reward function is equal to the sum of the Euclidean distance
625 between the objects and their goal positions and the distance between the arm and its goal position.