

Adversarial Intrinsic Motivation for Reinforcement Learning

Ishan Durugkar

Department of Computer Science
The University of Texas at Austin
Austin, TX, USA 78703
ishand@cs.utexas.edu

Mauricio Tec

Department of Statistics and Data Sciences
The University of Texas at Austin
Austin, TX, USA 78703
mauriciogtec@utexas.edu

Scott Niekum

Department of Computer Science
The University of Texas at Austin
Austin, TX, USA 78703
sniekum@cs.utexas.edu

Peter Stone

Department of Computer Science
The University of Texas at Austin
Austin, TX, USA 78703 and
Sony AI
pstone@cs.utexas.edu

Abstract

Learning with an objective to minimize the mismatch with a reference distribution has been shown to be useful for generative modeling and imitation learning. In this paper, we investigate whether one such objective, the Wasserstein-1 distance between a policy’s state visitation distribution and a target distribution, can be utilized effectively for reinforcement learning (RL) tasks. Specifically, this paper focuses on goal-conditioned reinforcement learning where the idealized (unachievable) target distribution has full measure at the goal. This paper introduces a quasimetric specific to Markov Decision Processes (MDPs) and uses this quasimetric to estimate the above Wasserstein-1 distance. It further shows that the policy that minimizes this Wasserstein-1 distance is the policy that reaches the goal in as few steps as possible. Our approach, termed Adversarial Intrinsic Motivation (AIM), estimates this Wasserstein-1 distance through its dual objective and uses it to compute a supplemental reward function. Our experiments show that this reward function changes smoothly with respect to transitions in the MDP and directs the agent’s exploration to find the goal efficiently. Additionally, we combine AIM with Hindsight Experience Replay (HER) and show that the resulting algorithm accelerates learning significantly on several simulated robotics tasks when compared to other rewards that encourage exploration or accelerate learning.

1 Introduction

Reinforcement Learning (RL) [70] deals with the problem of learning a policy to accomplish a given task in an optimal manner. This task is typically communicated to the agent by means of a reward function. If the reward function is sparse [4] (e.g., most transitions yield a reward of 0), much random exploration might be needed before the agent experiences any signal relevant to learning [11, 2].

Some of the different ways to speed up reinforcement learning by modifying or augmenting the reward function are shaped rewards [51], redistributed rewards [2], intrinsic motivations [8, 65, 67, 68, 53, 56], and learned rewards [77, 53]. Unfortunately, the optimal policy under such modified rewards might sometimes be different than the optimal policy under the task reward [51, 18]. The problem of

learning a reward signal that speeds up learning by communicating *what to do* but does not interfere by specifying *how to do it* is thus a useful and complex one [78].

This work considers whether a task-dependent reward function learned based on the distribution mismatch between the agent’s state visitation distribution and a target task (expressed as a distribution) can guide the agent towards accomplishing this task. Adversarial methods to minimize distribution mismatch have been used with great success in generative modeling [28] and imitation learning [38, 24, 75, 72, 27]. In both these scenarios, the task is generally to minimize the mismatch with a target distribution induced by the data or expert demonstrations. Instead, we consider the task of goal-conditioned RL, where the ideal target distribution assigns full measure to a goal state. While the agent can never match this idealized target distribution perfectly unless starting at the goal, intuitively, minimizing the mismatch with this distribution should lead to trajectories that maximize the proportion of the time spent at the goal, thereby prioritizing transitions essential to doing so.

The theory of optimal transport [74] gives us a way to measure the distance between two distributions (called the Wasserstein distance) even if they have disjoint support. Previous work [3, 31] has shown how a neural network approximating a potential function may be used to estimate the Wasserstein-1 distance using its dual formulation, but assumes that the metric space this distance is calculated on is Euclidean. A Euclidean metric might not be the appropriate metric to use in more general RL tasks however, such as navigating in a maze or environments where the state features change sharply with transitions in the environment.

This paper introduces a quasimetric tailored to Markov Decision Processes (MDPs), the time-step metric, to measure the Wasserstein distance between the agent’s state visitation distribution and the idealized target distribution. While this time-step metric could be an informative reward on its own, estimating it is a problem as hard as policy evaluation [30]. Instead, we show that the dual objective, which maximizes difference in potentials while utilizing the structure of this quasimetric for the necessary regularization, can be optimized through sampled transitions.

We use this dual objective to estimate the Wasserstein-1 distance and propose a reward function based on this estimated distance. An agent that maximizes returns under this reward minimizes this Wasserstein-1 distance. The competing objectives of maximizing the difference in potentials for estimating the Wasserstein distance and minimizing it through reinforcement learning on the subsequent reward function leads to our algorithm, Adversarial Intrinsic Motivation (AIM).

Our analysis shows that if the above Wasserstein-1 distance is computed using the time-step metric, then minimizing it leads to a policy that reaches the goal in the minimum expected number of steps. It also shows that if the environment dynamics are deterministic, then this policy is the optimal policy.

In practice, minimizing the Wasserstein distance works well even when the environment dynamics are stochastic. Our experiments show that AIM learns a reward function that changes smoothly with transitions in the environment. We further conduct experiments on the family of goal-conditioned reinforcement learning problems [1, 61] and show that AIM when used along with hindsight experience replay (HER) greatly accelerates learning of an effective goal-conditioned policy compared to learning with HER and the sparse task reward. Further, our experiments show that this acceleration is similar to the acceleration observed by using the actual distance to the goal as a dense reward.

2 Related Work

We highlight the related work based on the various aspects of learning that this work touches, namely intrinsic motivation, goal-conditioned reinforcement learning, and adversarial imitation learning.

2.1 Intrinsic Motivation

Intrinsic motivations [8, 56, 55] are rewards presented by an agent to itself in addition to the external task-specific reward. Researchers have pointed out that such intrinsic motivations are a characteristic of naturally intelligent and curious agents [29, 5, 6]. Intrinsic motivation has been proposed as a way to encourage RL agents to learn skills [10, 9, 64, 60] that might be useful across a variety of tasks, or as a way to encourage exploration [11, 63, 7, 23]. The optimal reward framework [65, 68] and shaped rewards [51] (if generated by the agent itself) also consider intrinsic motivation as a way to assist an RL agent in learning the optimal policy for a given task. Such an intrinsically motivated

reward signal has previously been learned through various methods such as evolutionary techniques [53, 62], meta-gradient approaches [67, 77, 78], and others. The Wasserstein distance has been used to present a valid reward for imitation learning [75, 19] as well as program synthesis [26].

2.2 Goal-Conditioned Reinforcement Learning

Goal-conditioned reinforcement learning [42] can be considered a form of multi-task reinforcement learning [17] where the agent is given the goal state it needs to reach at the beginning of every episode, and the reward function is sparse with a non-zero reward only on reaching the goal state. UVFA [61], HER [1], and others [76, 20] consider this problem of reaching certain states in the environment. Relevant to our work, Venkattaramanujam et al. [73] learns a distance between states using a random walk that is then used to shape rewards and speed up learning, but requires goals to be visited before the distance estimate is useful. DisCo RL [50] extends the idea of goal-conditioned RL to distribution-conditioned RL.

Contemporaneously, Eysenbach et al. [21, 22] has proposed a method which considers goals and examples of success and tries to predict and maximize the likelihood of seeing those examples under the current policy and trajectory. For successful training, this approach needs the agent to actually experience the goals or successes. Their solution minimizes the Hellinger distance to the goal, a form of f -divergence. AIM instead uses the Wasserstein distance which is theoretically more informative when considering distributions that are disjoint, and does not require the assumption that the agent has already reached the goal through random exploration. Our experiments in fact verify the hypothesis that AIM induces a form of directed exploration in order to reach the goal.

2.3 Adversarial Imitation Learning and Minimizing Distribution Mismatch

Adversarial imitation learning [38, 24, 75, 72, 27] has been shown to be an effective method to learn agent policies that minimize distribution mismatch between an agent’s state-action visitation distribution and the state-action visitation distribution induced by an expert’s trajectories. In most cases this distribution that the expert induces is achievable by the agent and hence these techniques aim to match the expert distribution exactly. In the context of goal-conditioned reinforcement learning, GoalGAIL [20] uses adversarial imitation learning with a few expert demonstrations to accelerate the learning of a goal-conditioned policy. In this work, we focus on unrealizable target distributions that cannot be completely matched by the agent, and indeed, are not induced by any trajectory distribution.

FAIRL [27] is an adversarial imitation learning technique which minimizes the Forward KL divergence and has been shown experimentally to cover some hand-specified state distributions, given a smoothness regularization as used by WGAN [31]. f -IRL [52] learns a reward function where the optimal policy matches the expert distribution under the more general family of f -divergences. Further, techniques beyond imitation learning [45, 36] have looked at matching a uniform distribution over states to guarantee efficient exploration.

3 Background

In this section we first set up the goal-conditioned reinforcement learning problem, and then give a brief overview of optimal transport.

3.1 Goal-Conditioned Reinforcement Learning

Consider a goal-conditioned MDP as the tuple $(S, A, G, P, \rho_0, \sigma, \gamma)$ with discrete state space S , discrete action space A , a subset of states which is the goal set $G \subseteq S$, and transition dynamics $P : S \times A \times G \rightarrow \Delta(S)$ ($\Delta(\cdot)$ is a distribution over a set) which might vary based on the goal (see below). $\rho_0 : \Delta(S)$ is the starting state distribution, and $\sigma : \Delta(G)$ is the distribution a goal is drawn from. $\gamma \in [0, 1)$ is the discount factor. We use discrete states and actions for ease of exposition, but our idea extends to continuous states and actions, as seen in the experiments.

At the beginning of an episode, the starting state is drawn from ρ_0 and the goal for that episode is drawn from σ . The reward function $r : S \times A \times S \times G \rightarrow \mathbb{R}$ is deterministic, and $r(s_t, a_t, s_{t+1} | s_g) := \mathbb{1}[s_{t+1} = s_g]$. That is, there is a positive reward when an agent reaches the goal ($s_{t+1} = s_g$), and 0 everywhere else. Since the goal is given to the agent at the beginning of

the episode, in goal-conditioned RL the agent knows what this task reward function is (unlike the more general RL problem). The transition dynamics are goal-conditioned as well, with an automatic transition to an absorbing state \bar{s} on reaching the goal s_g and then staying in that state with no rewards thereafter ($P(\bar{s}|s_g, a, s_g) = 1 \forall a \in \mathcal{A}$ and $P(\bar{s}|\bar{s}, a, s_g) = 1 \forall a \in \mathcal{A}$). In short, the episode terminates on reaching the goal state.

The agent takes actions in this environment based on a policy $\pi \in \Pi : \mathcal{S} \times \mathcal{G} \rightarrow \Delta(\mathcal{A})$. The return H_g for an episode with goal s_g is the discounted cumulative reward over that episode $H_g = \sum_{t=0}^T \gamma^t r(s_t, a_t, s_{t+1}|s_g)$, where $s_0 \sim \rho_0$, $a_t \sim \pi(j_{s_t}, s_g)$, and $s_{t+1} \sim P(j_{s_t}, a_t, s_g)$. The agent aims to find the policy $\pi = \arg \max_{\pi \in \Pi} \mathbb{E}_{g \sim \mathcal{G}} \mathbb{E}_{s_0 \sim \rho_0} \mathbb{E}_{\pi} [H_g]$ that maximizes the expected returns in this goal-conditioned MDP. For a policy π , the agent’s goal-conditioned state distribution $\rho_{\pi}(s|s_g) = \mathbb{E}_{s_0 \sim \rho_0} [(1 - \gamma) \sum_{t=0}^T \gamma^t P(s_t = s | \pi, s_g)]$. Overloading the terminology a bit, we also define the goal-conditioned target distribution $\rho_g(s|s_g) = \delta(s_g)$, a Dirac measure at the goal state s_g .

While learning using traditional RL paradigms is possible in goal-conditioned RL, there has also been previous work (Section 2.2) on leveraging the structure of the problem across goals. Hindsight Experience Replay (HER) [1] attempts to speed up learning in this sparse reward setting by taking episodes of agent interactions, where they might not have reached the goal specified for that episode, and relabeling the transitions with the goals that *were* achieved during the episode. Off-policy learning algorithms are then used to learn from this relabeled experience.

3.2 Optimal Transport and Wasserstein-1 Distance

The theory of optimal transport [74, 14] considers the question of how much work must be done to transport one distribution to another optimally, where this notion of work is defined by the use of a ground metric d . More concretely, consider a metric space (M, d) where M is a set and d is a metric on M (Definitions in Appendix A). For two distributions μ and ν with finite moments on the set M , the Wasserstein- p distance is denoted by:

$$W_p(\mu, \nu) := \inf_{\zeta \in \mathcal{Z}(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \zeta} [d(X, Y)^p]^{1/p} \quad (1)$$

where \mathcal{Z} is the space of all possible couplings, i.e. joint distributions $\zeta \in \Delta(M \times M)$ whose marginals are μ and ν respectively. Finding this optimal coupling tells us what is the least amount of work, as measured by d , that needs to be done to convert μ to ν . This Wasserstein- p distance can then be used as a cost function (negative reward) by an RL agent to match a given target distribution [75, 19, 26].

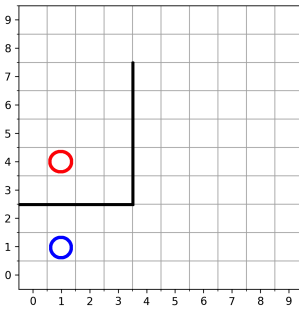


Figure 1: Grid world example

Finding the ideal coupling above is generally considered intractable. However, if what we need is an accurate estimate of the Wasserstein distance and not the optimal transport plan we can turn our attention to the dual form of the Wasserstein-1 distance. The Kantorovich-Rubinstein duality [74, 57] for the Wasserstein-1 distance (which we refer to simply as the Wasserstein distance hereafter) on a ground metric d gives us:

$$W_1(\mu, \nu) = \sup_{\text{Lip}(f) \leq 1} \mathbb{E}_{y \sim \nu} [f(y)] - \mathbb{E}_{x \sim \mu} [f(x)] \quad (2)$$

where the supremum is over all 1-Lipschitz functions $f : M \rightarrow \mathbb{R}$ in the metric space. Importantly, Jevtić [40] has recently shown that this dual formulation extends to quasimetric spaces as well.

More details, such as the definition of the Lipschitz constant of a function and special cases used in WGAN [3, 31, 27] are elaborated in Appendix B. One last note of importance is that the Lipschitz constant of the potential function f is computed based on the ground metric d .

4 Time-Step Metric

The choice of the ground metric d is important when computing the Wasserstein distance between two distributions. That is, if we want the Wasserstein distance to give an estimate of the work needed to transport the agent’s state visitation distribution to the goal state, the ground metric should incorporate a notion of this work.

Consider the grid world shown in Figure 1, where a wall (bold line) marks an impassable barrier in part of the state space. If the states are specified by their Cartesian coordinates on the grid, the Manhattan distance between the states specified by the blue and red circles is not representative of the optimal cost to go from one to the other. This mismatch would lead to an underestimation of the work involved if the two distributions compared were concentrated at those two circles. Similarly, there will be errors in estimating the Wasserstein distance if the grid world is toroidal (where an agent is transported to the opposite side of the grid if it walks off one side) or if the transitions are asymmetric (windy grid world [70]).

To estimate the work needed to transport measure in an MDP when executing a policy π , we consider a *quasimetric* – a metric that does not need to be symmetric – dependent on the number of transitions experienced before reaching the goal when executing that policy.

Definition 1. The *time-step metric* d_T^π in an MDP with state space S , action space A , transition function P , and agent policy π is a quasimetric where the distance from state $s \in S$ to state $s_g \in S$ is based on the expected number of transitions under policy π .

$$d_T^\pi(s, s_g) := \mathbb{E} [T(s_g/\pi, s)]$$

where $T(s_g/\pi, s)$ is the random variable for the first time-step that state s_g is encountered by the agent after starting in state s and following policy π .

This quasimetric has the property that the per step cost is uniformly 1 for all transitions except ones from the goal to the absorbing state (and the absorbing state to itself), which are 0. Thus, it can be written recursively as:

$$d_T^\pi(s, s_g) = \begin{cases} 0 & \text{if } s = s_g \\ 1 + \mathbb{E}_{a \sim \pi(\cdot|s, s_g)} \mathbb{E}_{s' \sim P(\cdot|s, a, s_g)} [d_T^\pi(s', s_g)] & \text{otherwise} \end{cases} \quad (3)$$

Recall that in order to estimate the Wasserstein distance using the dual (Equation 2) in a metric space where the ground metric d is this time-step metric, the potential function $f : S \rightarrow \mathbb{R}$ needs to be 1-Lipschitz with respect to d_T^π . In Appendix C we prove that L -Lipschitz continuity can be ensured by enforcing that the difference in values of f on expected transitions from every state are bounded by L , implying

$$\text{Lip}(f) \leq \sup_{s \in S} \left\{ \mathbb{E}_{a \sim \pi(\cdot|s, s_g)} \mathbb{E}_{s' \sim P(\cdot|s, a, s_g)} [f(s') - f(s)] \right\}. \quad (4)$$

Note that finding a proper way to enforce the Lipschitz constraint in adversarial methods remains an open problem [48]. However, for the time-step metric considered here, equation 4 is one elegant way of doing so. By ensuring that the Kantorovich potentials do not drift too far from each other on expected transitions under agent policy π in the MDP, the conditions necessary for the potential function to estimate the Wasserstein distance can be maintained [74, 3]. Finally, the minimum distance d_T from state s to a given goal state s_g (corresponding to policy π) is defined by the Bellman optimality condition (Equation 16 in Appendix D).

Consider how the time-step distance to the goal and the value function for goal-conditioned RL relate to each other. When the reward is 0 everywhere except for transitions to the goal state, the value becomes $V^\pi(s/s_g) = \mathbb{E} [\gamma^{T(s_g/\pi, s)}]$. $d_T^\pi(s_0, s_g)$ and $V(s_0/s_g)$ are related as follows.

Proposition 1. A lower bound on the value of any state under a policy π can be expressed in terms of the time-step distance from that state to the goal: $V(s_0/s_g) \geq \gamma^{d_T^\pi(s_0, s_g)}$.

The proofs for all theoretical results are in Appendix D. The Jensen gap $\Delta_{\text{Jensen}}^\pi(s) := V^\pi(s/s_g) - \gamma^{d_T^\pi(s, s_g)}$ describes the sharpness of the lower bound in the proposition above and it is zero if and only if $\text{Var}(T(s_g/\pi, s)) = 0$ [46]. From this line of reasoning, we deduce the following proposition:

Proposition 2. If the transition dynamics are deterministic, the policy that maximizes expected return is the policy that minimizes the time-step metric ($\pi^* = \pi^*$).

5 Wasserstein-1 Distance for Goal-Conditioned Reinforcement Learning

In this section we consider the problem of goal-conditioned reinforcement learning. In Section 5.1 we analyze the Wasserstein distance computed under the time-step metric d_T^π . Section 5.2 proposes an

algorithm, Adversarial Intrinsic Motivation (AIM), to learn the potential function for the Kantorovich-Rubinstein duality used to estimate the Wasserstein distance, and giving an intrinsic reward function used to update the agent policy in tandem.

5.1 Wasserstein-1 Distance under the Time-Step Metric

From Sections 3.2 and 4 the Wasserstein distance under the time-step metric d_T^π of an agent policy π with visitation measure ρ_π to a particular goal s_g and its distribution ρ_g can be expressed as:

$$W_1^\pi(\rho_\pi, \rho_g) = \sum_{s \in \mathcal{S}} \rho_\pi(s) d_T^\pi(s, s_g) \quad (5)$$

where W_1^π refers to the Wasserstein distance with the ground metric d_T^π .

The following proposition shows that the Wasserstein distance decreases as $d_T^\pi(s, s_g)$ decreases, while also revealing a surprising connection with the Jensen gap.

Proposition 3. *For a given policy π , the Wasserstein distance of the state visitation measure of that policy from the goal state distribution ρ_g under the ground metric d_T^π can be written as*

$$W_1^\pi(\rho_\pi, \rho_g) = \mathbb{E}_{s_0 \sim \rho_0} \left[h(d_T^\pi(s_0, s_g)) + \frac{\gamma}{1-\gamma} (\Delta_{Jensen}^\pi(s_0) - 1) \right] \quad (6)$$

where h is an increasing function of d_T^π .

The first component in the above analytical expression shows that the Wasserstein distance depends on the expected number of steps, decreasing if the expected distance decreases. The second component shows the risk-averse nature of the Wasserstein distance. Concretely, the bounds for the Jensen inequality given by Liao and Berg [46] imply that there are non-negative constants $C_1 = C_1(d_T^\pi, \gamma)$ and $C_2 = C_2(d_T^\pi, \gamma)$ depending only on the expected distance and discount factor such that

$$C_1 \text{Var}(T(s_g | \pi, s)) \leq \Delta_{Jensen}^\pi(s) \leq C_2 \text{Var}(T(s_g | \pi, s)).$$

From the above, we can deduce that a policy with lower variance will have lower Wasserstein distance when compared to a policy with same expected distance from the start but higher variance. The relation between the optimal policy in goal-conditioned RL and the Wasserstein distance can be made concrete if we consider deterministic dynamics.

Theorem 1. *If the transition dynamics are deterministic, the policy that minimizes the Wasserstein distance over the time-step metrics in a goal-conditioned MDP (see equation 5) is the optimal policy.*

5.2 Adversarial Intrinsic Motivation to minimize Wasserstein-1 Distance

The above section makes it clear that minimizing the Wasserstein distance to the goal will lead to a policy that reaches the goal in as few steps as possible in expectation. If the dynamics of the MDP are deterministic, this policy will also be optimal. Note that the dual form (Equation 2) can be used to estimate the distance, *even if the ground metric d_T^π is not known*. The smoothness requirement on the potential function f can be ensured with the constraint in Equation 4 on all states and subsequent transitions expected under the agent policy.

Now consider the full problem. The reinforcement learning algorithm aims to learn a goal-conditioned policy with parameters $\theta \in \Theta$ whose state visitation distribution ρ_θ minimizes the Wasserstein distance to a goal-conditioned target distribution ρ_g for a given goal $s_g \in \mathcal{S}$. AIM leverages the presence of the set of goals that the agent should be capable of reaching with a goal-conditioned potential function $f_\phi : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ with parameters $\phi \in \Phi$. These objectives of the potential function and the agent can be expressed together using the following adversarial objective:

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathbb{E}_{s_g \sim \rho_g} \left[f_\phi(s_g, s_g) - \mathbb{E}_{s \sim \rho_\theta} [f_\phi(s, s_g)] \right] \quad (7)$$

where the potential function f_ϕ is 1-Lipschitz over the state space. Combining the objectives in Equations 7 and 4, the loss for the potential function f_ϕ then becomes:

$$L_f := \mathbb{E}_{s_g \sim \rho_g} \left[f_\phi(s_g, s_g) + \mathbb{E}_{s \sim \rho_\theta} [f_\phi(s, s_g)] \right] + \lambda \mathbb{E}_{(s, a, s', s_g) \sim \mathcal{D}} \left[(\max(|f_\phi(s, s_g) - f_\phi(s', s_g)|, 1), 0)^2 \right] \quad (8)$$

Where the distribution D should ideally contain all states in S , expected goals in G , and the transitions according to the agent policy π_θ and transition function P . Such a distribution is difficult to obtain directly. AIM approximates it with a small replay buffer of transitions from recent episodes experienced by the agent, and relabels these episodes with achieved goals (similar to HER [1]). Such an approximation does not respect the discounted measure of states later on in an episode, but is consistent with how other approaches in deep reinforcement learning tend to approximate the state visitation distribution, especially for policy gradient approaches [54]. While it does not include all states and all goals, we see empirically that the above approximation works well.

Now we turn to the reward function that should be presented to the agent such that maximizing the return will minimize the Wasserstein distance. The Wasserstein discriminator is a potential function [51] (its value depends on the state). It can thus be used to create a shaped reward $\hat{r}(s, a, s^o, s_g) = r(s, a, s^o, s_g) + \gamma f_\phi(s^o, s_g) - f_\phi(s, s_g)$ without risk of changing the optimal policy. Alternatively, we can explicitly minimize samples of the Wasserstein distance: $\hat{r}(s, a, s^o, s_g) = f_\phi(s^o, s_g) - f_\phi(s, s_g)$. Finally, instead of the second term $f_\phi(s, s_g)$, we can just use a constant bias term. In practice, all these choices work well, and the experiments use the latter (with $b = \max_{s \in S} f_\phi(s, s_g)$) to reduce variance in \hat{r} .

$$\hat{r}(s, a, s^o, s_g) = f_\phi(s^o, s_g) - b \quad (9)$$

The basic procedure to learn and use adversarial intrinsic motivation (AIM) is laid out in Algorithm 1, and also includes how to use this algorithm in conjunction with HER. If not using HER, Line 23 where hindsight goals are added to the replay buffer can be skipped.

6 Experiments

Our experiments evaluate the extent to which the reward learned through AIM is useful as a proxy for the environment reward signal, or in tandem with the environment reward signal. In particular, we ask the following questions:

- Does AIM speed up learning of a policy to get to a single goal compared to learning with a sparse reward?
- Does the learned reward function qualitatively guide the agent to the goal?
- Does AIM work well with stochastic transition dynamics or sharp changes in the state features?
- Does AIM generalize to a large set of goals and continuous state and action spaces?

Our experiments suggest that the answer to all 4 questions is “yes”, with the first three questions tested in the grid world presented in Figure 1 where the goal is within a room, and the agent has to go around the room from its start state to reach the goal. Goal-conditioned tasks in the established Fetch robot domain show that AIM also accelerates learning across multiple goals in continuous state and action spaces.

This section compares an agent learning with a reward learned through AIM with other intrinsic motivation signals that induce general exploration or shaped rewards that try to guide the agent to the goal. The experiments show that AIM guides the agent’s exploration more efficiently and effectively than a general exploration bonus, and adapts to the dynamics of the environment better than other techniques we compare to. As an overview, the baselines we compare to are:

- RND**: with random network distillation (RND) [16] used to provide a general exploration bonus.
- MC**: with the distance between states learned through regression of Monte Carlo rollouts of the agent policy, similar to Hartikainen et al. [34].
- SMiRL**: SMiRL [13] is used to provide a bonus intrinsic motivation reward that minimizes the overall surprise in an episode.
- DiscoRL** The DiscoRL [50] approach presents a reward to maximize the likelihood of a target distribution (normal distribution at the goal). In practice this approach is equivalent to a negative L2 distance to the goal, which we compare to in the grid world domain.
- GAIL**: additional GAIL [38] rewards using trajectories relabeled with achieved goals considered as having come from the expert in hindsight. This baseline is compared to in the Fetch robot domain, since that is the domain where we utilize hindsight relabeling.

Grid World In this task, the goal is inside a room and the agent’s starting position is such that it needs to navigate around the room to find the doorway and be able to reach the goal. The agent can

move in the 4 cardinal directions unless blocked by a wall or the edge of the grid. The agent policy is learned using soft Q-learning [32] with no hindsight goals used for this experiment.

The agent’s state visitation distribution after just 100 Q-function updates when using AIM-learned rewards is shown in Figure 2a and the learned rewards for each state are plotted in Figure 2b. The state visitation when learning with the true task reward shows that the agent is unable to learn a policy to the goal (Figure 2c). These figures show that AIM enables the agent to reach the goal and learn the required policy quickly, while learning with the sparse task reward fails to do so.

In Appendix F we also compare to the baselines described above and show that AIM learns a reward that is more efficient at directing the agent’s exploration and more flexible to variations of the environment dynamics, such as stochastic dynamics or transitions that cause a sharp change in the state features. None of the baselines compared to were able to direct the agent to the goal in this grid world even after given up to 5 more interactions with the environment to train. AIM’s use of the time-step metric also enabled it to adapt to variations of the environment dynamics better than the gradient penalty based regularization used in Wasserstein GANs [31] and adversarial imitation learning [27] approaches.

Fetch Robot The generalization capability of AIM across multiple goals in goal-conditioned RL tasks with continuous states and actions is tested in the MuJoCo simulator [71], on the Fetch robot tasks from OpenAI gym [15] which have been used to evaluate learning of goal-conditioned policies previously [1, 76]. Descriptions of these tasks and their goal space is in Appendix H. We soften the Dirac target distribution for continuous states to instead be a Gaussian with variance of 0.01 of the range of each feature.

The goals in this setting are not the full state, but rather the dimensions of factored states relevant to the given goal. The task wrapper additionally returns the features of the agent’s state in this reduced goal space, and so AIM can use it to learn our reward function, rather than the full state space. It is unclear how this smaller goal space might affect AIM. While the smaller goal space might make learning easier for potential function f_ϕ , the partially observable nature of the goals might lead to a less informative reward.

We combine AIM with HER (refer subsection 3.1) and refer to it as [AIM + HER]. We compare this agent to the baselines we referred to above, as well as the sparse environment reward (R + HER) and the dense reward derived from the negative Euclidean (L_2) distance to the goal ($-L_2$ + HER). The L_2 distance is proportional to the number of steps it should take the agent to reach the goal in this environment, and so the reward based on it can act as an oracle reward that we can use to test how efficiently AIM can learn a reward function that helps the agent learn its policy. We used the HER implementation using Twin Delayed DDPG (TD3) [25] as the underlying RL algorithm from the stable baselines repository [37]. We did an extensive sweep of the hyperparameters for the baseline HER + R (laid out in Appendix H), with a coarser search on relevant hyperparameters for AIM.

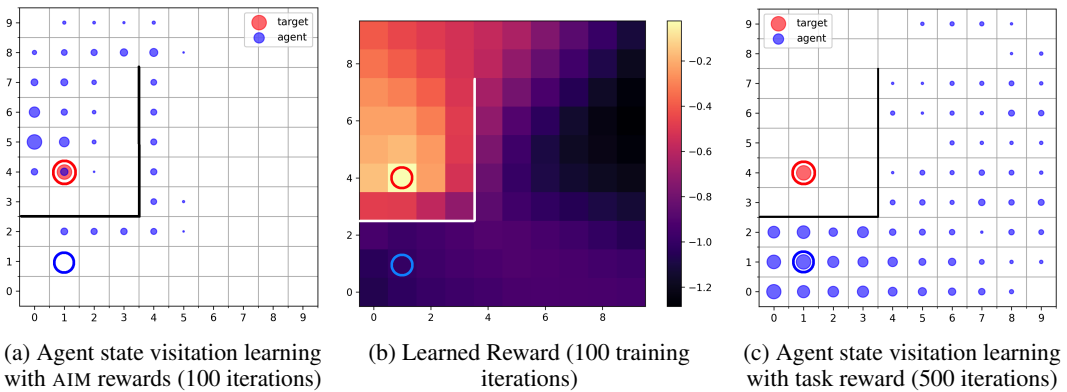


Figure 2: Grid world experiments. Agent’s undiscounted state visitation (2a, 2c): Blue circle indicates agent’s start state. Red circle is the goal. Blue bubbles indicate relative time agent’s policy causes it to spend in respective states. Learned reward function (2b): AIM reward at each state of the grid world. Bold black (or white) lines indicate walls the agent cannot transition through.

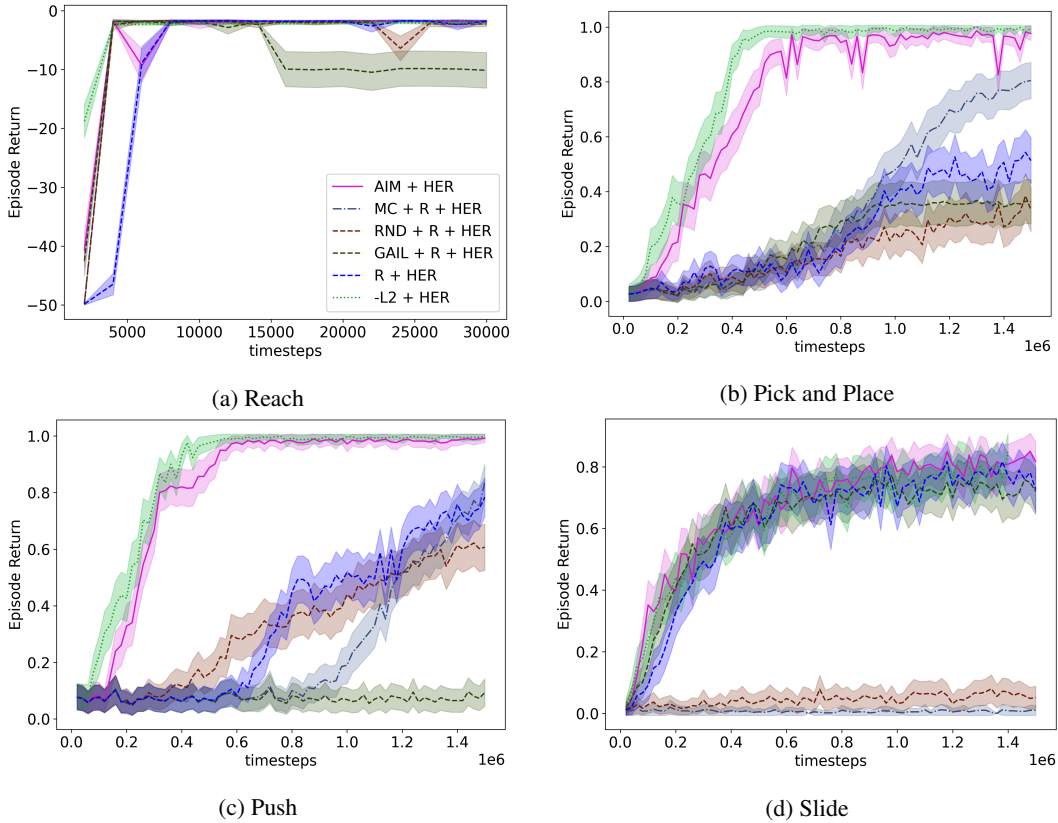


Figure 3: Evaluating AIM with HER on some goal-conditioned RL tasks in the Fetch domain. AIM learns the reward function in tandem with the policy updates. The “ L_2 ” reward is the true negative distance to goal, acting as an oracle reward in this domain. The other baselines are detailed above.

Figure 3 shows that using the AIM-learned reward speeds up learning in three of the four Fetch domains, even without the environment reward. This improvement is very close to what we would see if we used the dense reward (based on the actual distance function in this domain). An additional comparison with an agent learning with both the AIM-learned reward as well as the task reward (AIM + R + HER) can be seen in Figure 7 in the Appendix, showing that using both signals accelerates learning even more. These results also highlight that AIM continues to work in continuous state and action spaces, even though our analysis focuses on discrete states and actions. Results are averaged across the 6 different seeds, with shaded regions showing standard error across runs. Statistical analysis using a mixed effects ANOVA and a Tukey test at a significance level of 95% (more detail in Appendix G) show that in three of the four environments AIM and AIM+ R have similar odds of reaching the goal as the dense shaped reward, and in all four environments AIM and AIM+ R have higher odds of reaching the goal compared to the sparse reward.

The other baselines compare well to AIM in the Fetch Reach domain (Figure 3a), but do not do as well on the other problems. In fact, none of the other baselines outperform the vanilla baseline [R + HER] in all the domains. The RND rewards help the agent to start learning faster in the Push domain (Figure 3c), but lead to worse performance in Pick and Place (Figure 3b). On the other hand, learning the distance function through MC regression helps in the Pick and Place domain, but slows down learning when dealing with Push. Most notably, both these approaches cause learning to fail in the Slide domain (Figure 3d), where the credit assignment problem is especially difficult. GAIL works as well as AIM and the vanilla baseline in Slide, but underperforms in the other domains. We hypothesize that the additional rewards in these baselines conflict with the task reward. Additionally, none of the three new baselines work well if we do not provide the task reward in addition to the specific bonus for that algorithm.

We did not find any configuration in the Fetch Reach domain where [SMiRL + R + HER] was able to accomplish the task in the given training budget. Since SMiRL did not work on the grid world or Fetch Reach, we did not try it out on any of the other domains.

FAIRL [27] (which has been shown to learn policies that cover hand-specified state distributions) was also applied on these 4 domains but it failed to learn at all. Interestingly, scaling the reward such that it is always negative led to similar performance to (but not better than) AIM. We hypothesize that FAIRL, as defined and presented, fails in these domains because the environments are episodic, and the episode ends earlier if the goal is reached. Since the FAIRL reward is positive closer to the target distribution, the agent can get close to the target, but refrain from reaching it (and ending the episode) to collect additional positive reward.

The domain where AIM does not seem to have a large advantage (Slide) is one where the agent strikes an object initially and that object has to come to rest near the goal. In fact, AIM-learned rewards, the vanilla environment reward R , and the oracle L_2 rewards all lead to similar learning behavior, indicating that this particular task does not benefit much from shaped rewards. The reason for this invariance might be that credit assignment has to propagate back to the single point when the agent strikes the object regardless of how dense the subsequent reward is.

7 Discussion and Future Work

Approaches for estimating the Wasserstein distance to a target distribution by considering the dual of the Kantorovich relaxation have been previously proposed [3, 31, 75], but assume that the ground metric is the L_2 distance. We improve upon them by choosing a metric space more suited to the MDP and notions of optimality in the MDP. This choice allows us to leverage the structure introduced by the dynamics of the MDP to regularize the Kantorovich potential using a novel objective.

Previous work [12] has pointed out that the gradients from sample estimates of the Wasserstein distance might be biased. This issue is mitigated in our implementation through multiple updates of the discriminator, which they found to be empirically useful in reducing the bias. Additionally, recent work has pointed out that the discriminator in WGAN might be bad at estimating the Wasserstein distance [69]. While our experiments indicate that the potential function in AIM is learned appropriately, future work could look more deeply to verify possible inefficiencies in this estimation.

The process of learning the Wasserstein distance through samples of the environment while simultaneously estimating the cost of the full path is reminiscent of the A algorithm [33], where the optimistic heuristic encourages the agent to explore in a directed manner, and adjusts its estimates based on these explorations.

The discriminator objective (Equation 8) also bears some resemblance to a linear program formulation of the RL problem [58]. The difference is that this formulation minimizes the value function on states visited by the agent, while AIM additionally maximizes the potential at the goal state. This crucial difference has two main consequences. First, the potential function during learning is not equivalent to the value of the agent’s policy (verified by using this potential as a critic). Second, increasing the potential of the goal state in AIM directs the agent exploration in a particular direction (namely, the direction of sharpest increase in potential).

In the goal-conditioned RL setting, AIM seems to be an effective intrinsic reward that balances exploration and exploitation for the task at hand. The next step is to consider whether the Wasserstein distance can be estimated similarly for more general tasks, and whether minimizing this distance in those tasks leads to the optimal policy. A different potential avenue for future work is the problem of more general exploration [36, 45] by specifying a uniform distribution as the target, or using this directed exploration as an intermediate step for efficient exploration [41].

Finally, reward design is an important aspect of practical reinforcement learning. Not only do properly shaped reward speed up learning [51], but reward design can also subtly influence the kinds of behaviors deemed acceptable for the RL agent [44] and could be a potential safety issue keeping reinforcement learning from being deployed on real world problems. Learning-based approaches that can assist in specifying reward functions safely given alternative approaches for communicating the task could be of value in such a process of reward design, and an avenue for future research.

Acknowledgements and Funding Information

We thank Caroline Wang, Garrett Warnell, and Elad Liebman for discussion and feedback on this work. We also thank the reviewers for their thoughtful comments and suggestions that have helped to improve this paper.

This work has taken place in part in the Learning Agents Research Group (LARG) at the Artificial Intelligence Laboratory, and in part in the Personal Autonomous Robotics Lab (PeARL) at The University of Texas at Austin. LARG research is supported in part by the National Science Foundation (CPS-1739964, IIS-1724157, FAIN-2019844), the Office of Naval Research (N00014-18-2243), Army Research Office (W911NF-19-2-0333), DARPA, Lockheed Martin, General Motors, Bosch, and Good Systems, a research grand challenge at the University of Texas at Austin. PeARL research is supported in part by the NSF (IIS-1724157, IIS-1638107, IIS-1749204, IIS-1925082), ONR (N00014-18-2243), AFOSR (FA9550-20-1-0077), and ARO (78372-CS). This research was also sponsored by the Army Research Office under Cooperative Agreement Number W911NF-19-2-0333. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein. Peter Stone serves as the Executive Director of Sony AI America and receives financial compensation for this work. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

References

- [1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience Replay. *arXiv:1707.01495 [cs]*, February 2018. URL <http://arxiv.org/abs/1707.01495>. arXiv: 1707.01495.
- [2] Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. In *Advances in Neural Information Processing Systems*, pages 13566–13577, 2019.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein Generative Adversarial Networks. In *International Conference on Machine Learning*, pages 214–223, July 2017. URL <http://proceedings.mlr.press/v70/arjovsky17a.html>. ISSN: 2640-3498 Section: Machine Learning.
- [4] Dilip Arumugam, P. Henderson, and P. Bacon. An information-theoretic perspective on credit assignment in reinforcement learning. *ArXiv*, abs/2103.06224, 2021.
- [5] Gianluca Baldassarre. What are intrinsic motivations? a biological perspective. In *2011 IEEE international conference on development and learning (ICDL)*, volume 2, pages 1–8. IEEE, 2011.
- [6] Gianluca Baldassarre, Tom Stafford, Marco Mirolli, Peter Redgrave, Richard M Ryan, and Andrew Barto. Intrinsic motivations and open-ended development in animals, humans, and robots: an overview. *Frontiers in psychology*, 5:985, 2014.
- [7] Adrien Baranes and Pierre-Yves Oudeyer. R-iac: Robust intrinsically motivated exploration and active learning. *IEEE Transactions on Autonomous Mental Development*, 1(3):155–169, 2009.
- [8] Andrew G. Barto. Intrinsic Motivation and Reinforcement Learning. In Gianluca Baldassarre and Marco Mirolli, editors, *Intrinsically Motivated Learning in Natural and Artificial Systems*, pages 17–47. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-3-642-32374-4 978-3-642-32375-1. doi: 10.1007/978-3-642-32375-1_2. URL http://link.springer.com/10.1007/978-3-642-32375-1_2.
- [9] Andrew G Barto and Ozgür Simsek. Intrinsic motivation for reinforcement learning systems. In *Proceedings of the Thirteenth Yale Workshop on Adaptive and Learning Systems*, pages 113–118, 2005.
- [10] Andrew G Barto, Satinder Singh, and Nuttapon Chentanez. Intrinsically motivated learning of hierarchical collections of skills. In *Proceedings of the 3rd International Conference on Development and Learning*, pages 112–19. Cambridge, MA, 2004.

- [11] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in neural information processing systems*, pages 1471–1479, 2016.
- [12] Marc G. Bellemare, Ivo Danihelka, Will Dabney, S. Mohamed, Balaji Lakshminarayanan, Stephan Hoyer, and R. Munos. The cramer distance as a solution to biased wasserstein gradients. *ArXiv*, abs/1705.10743, 2017. URL <https://arxiv.org/abs/1705.10743>.
- [13] Glen Berseth, Daniel Geng, Coline Devin, Nicholas Rhinehart, Chelsea Finn, Dinesh Jayaraman, and Sergey Levine. Smirl: Surprise minimizing reinforcement learning in unstable environments. *arXiv preprint arXiv:1912.05510*, 2019.
- [14] Olivier Bousquet, Sylvain Gelly, Ilya Tolstikhin, Carl-Johann Simon-Gabriel, and Bernhard Schoelkopf. From optimal transport to generative modeling: the VEGAN cookbook. *arXiv:1705.07642 [stat]*, May 2017. URL <http://arxiv.org/abs/1705.07642>. arXiv: 1705.07642.
- [15] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [16] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.
- [17] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [18] Jack Clark and Dario Amodei. Faulty reward functions in the wild, Dec 2016. URL <https://openai.com/blog/faulty-reward-functions/>.
- [19] Robert Dadashi, Léonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal Wasserstein Imitation Learning. *arXiv:2006.04678 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2006.04678>. arXiv: 2006.04678.
- [20] Y. Ding, Carlos Florensa, Mariano Phielipp, and P. Abbeel. Goal-conditioned imitation learning. In *NeurIPS*, 2019.
- [21] Benjamin Eysenbach, Sergey Levine, and Ruslan Salakhutdinov. Replacing rewards with examples: Example-based policy search via recursive classification. *arXiv preprint arXiv:2103.12656*, 2021.
- [22] Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. C-learning: Learning to achieve goals via recursive classification. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=tc5qisoB-C>.
- [23] Sébastien Forestier, Rémy Portelas, Yoan Mollard, and Pierre-Yves Oudeyer. Intrinsically motivated goal exploration processes with automatic curriculum learning. *arXiv preprint arXiv:1708.02190*, 2017.
- [24] Justin Fu, Katie Luo, and Sergey Levine. Learning Robust Rewards with Adversarial Inverse Reinforcement Learning. *arXiv:1710.11248 [cs]*, August 2018. URL <http://arxiv.org/abs/1710.11248>. arXiv: 1710.11248.
- [25] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.
- [26] Yaroslav Ganin, Tejas Kulkarni, Igor Babuschkin, SM Ali Eslami, and Oriol Vinyals. Synthesizing programs for images using reinforced adversarial learning. In *International Conference on Machine Learning*, pages 1666–1675. PMLR, 2018.
- [27] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. In *Conference on Robot Learning*, pages 1259–1277. PMLR, 2020.
- [28] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- [29] Jacqueline Gottlieb, Pierre-Yves Oudeyer, Manuel Lopes, and Adrien Baranes. Information-seeking, curiosity, and attention: computational and neural mechanisms. *Trends in cognitive sciences*, 17(11):585–593, 2013.

- [30] Matthieu Guillot and Gautier Stauffer. The stochastic shortest path problem: a polyhedral combinatorics perspective. *European Journal of Operational Research*, 285(1):148–158, 2020.
- [31] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved Training of Wasserstein GANs. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5767–5777. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7159-improved-training-of-wasserstein-gans.pdf>.
- [32] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1352–1361, 2017.
- [33] Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE transactions on Systems Science and Cybernetics*, 4(2):100–107, 1968.
- [34] Kristian Hartikainen, Xinyang Geng, Tuomas Haarnoja, and Sergey Levine. Dynamical distance learning for semi-supervised and unsupervised skill discovery. In *International Conference on Learning Representations*, 2020.
- [35] Matthew Hausknecht and Peter Stone. Deep reinforcement learning in parameterized action space. In *International Conference on Learning Representations*, 2016. URL <https://arxiv.org/abs/1511.04143>.
- [36] Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy exploration. In *International Conference on Machine Learning*, pages 2681–2691. PMLR, 2019.
- [37] Ashley Hill, Antonin Raffin, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, Rene Traore, Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, and Yuhuai Wu. Stable baselines. <https://github.com/hill-a/stable-baselines>, 2018.
- [38] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29:4565–4573, 2016.
- [39] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010. ISSN 1533-7928. URL <http://jmlr.org/papers/v11/jaksch10a.html>.
- [40] Filip Jevtić. *Combinatorial Structure of Finite Metric Spaces*. PhD thesis, The University of Texas at Dallas, August 2018.
- [41] Yuu Jinnai, Jee Won Park, Marlos C. Machado, and George Konidaris. Exploration in reinforcement learning with deep covering options. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SkeIyaVtwB>.
- [42] Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, pages 1094–1099. Citeseer, 1993.
- [43] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2):209–232, 2002.
- [44] W Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward (mis) design for autonomous driving. *arXiv preprint arXiv:2104.13906*, 2021.
- [45] Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Ruslan Salakhutdinov. Efficient exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.
- [46] JG Liao and Arthur Berg. Sharpening jensen’s inequality. *The American Statistician*, 2018.
- [47] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016.
- [48] Kanglin Liu and Guoping Qiu. Lipschitz constrained gans via boundedness and continuity. *Neural Computing and Applications*, pages 1–13, 2020.
- [49] Douglas C Montgomery. *Design and analysis of experiments*. John wiley & sons, 2017.

- [50] Soroush Nasiriany. Disco rl: Distribution-conditioned reinforcement learning for general-purpose policies. Master’s thesis, EECS Department, University of California, Berkeley, Aug 2020. URL <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2020/EECS-2020-151.html>.
- [51] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287, 1999.
- [52] Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Benjamin Eysenbach. F-irl: Inverse reinforcement learning via state marginal matching. In *Conference on Robot Learning*, 2020.
- [53] Scott Niekum. Evolved intrinsic reward functions for reinforcement learning. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, pages 1955–1956, 2010.
- [54] Chris Nota and Philip S Thomas. Is the policy gradient a gradient? In *International Conference on Autonomous Agents and Multi-Agent Systems*, pages 939–947, 2020.
- [55] Pierre-Yves Oudeyer and Frederic Kaplan. How can we define intrinsic motivation? In *the 8th International Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic Systems*. Lund University Cognitive Studies, Lund: LUCS, Brighton, 2008.
- [56] Pierre-Yves Oudeyer and Frederic Kaplan. What is intrinsic motivation? a typology of computational approaches. *Frontiers in neurobotics*, 1:6, 2009.
- [57] Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019. URL <http://arxiv.org/abs/1803.00567>.
- [58] Martin L Puterman. Markov decision processes. *Handbooks in operations research and management science*, 2:331–434, 1990.
- [59] Antonin Raffin. RL baselines zoo. <https://github.com/araffin/rl-baselines-zoo>, 2018.
- [60] Vieri Giuliano Santucci, Gianluca Baldassarre, and Marco Mirolli. Which is the best intrinsic motivation signal for learning multiple skills? *Frontiers in neurobotics*, 7:22, 2013.
- [61] T. Schaul, Daniel Horgan, K. Gregor, and D. Silver. Universal value function approximators. In *ICML*, 2015.
- [62] Massimiliano Schembri, Marco Mirolli, and Gianluca Baldassarre. Evolving internal reinforcers for an intrinsically motivated reinforcement-learning robot. In *2007 IEEE 6th International Conference on Development and Learning*, pages 282–287. IEEE, 2007.
- [63] Özgür Şimşek and Andrew G Barto. An intrinsic reward mechanism for efficient exploration. In *Proceedings of the 23rd international conference on Machine learning*, pages 833–840, 2006.
- [64] Satinder Singh, Andrew G Barto, and Nuttapon Chentanez. Intrinsically motivated reinforcement learning. In *Advances in neural information processing systems*, pages 1281–1288, 2005.
- [65] Satinder Singh, Richard L. Lewis, Andrew G. Barto, and Jonathan Sorg. Intrinsically Motivated Reinforcement Learning: An Evolutionary Perspective. *IEEE Transactions on Autonomous Mental Development*, 2(2):70–82, June 2010. ISSN 1943-0612. doi: 10.1109/TAMD.2010.2051031. Conference Name: IEEE Transactions on Autonomous Mental Development.
- [66] M. Smyth. Quasi uniformities: Reconciling domains with metric spaces. In *MFPS*, 1987.
- [67] Jonathan Sorg, Richard L Lewis, and Satinder P Singh. Reward design via online gradient ascent. In *Advances in Neural Information Processing Systems*, pages 2190–2198, 2010.
- [68] Jonathan Sorg, Satinder P Singh, and Richard L Lewis. Internal rewards mitigate agent boundedness. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 1007–1014, 2010.
- [69] Jan Stanczuk, Christian Etmann, L. Kreuzer, and C. Schönlieb. Wasserstein gans work because they fail (to approximate the wasserstein distance). *ArXiv*, abs/2103.01678, 2021.
- [70] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- [71] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- [72] Faraz Torabi, Garrett Warnell, and Peter Stone. Generative Adversarial Imitation from Observation. *arXiv:1807.06158 [cs, stat]*, June 2019. URL <http://arxiv.org/abs/1807.06158>. arXiv: 1807.06158.
- [73] Srinivas Venkattaramanujam, E. Crawford, T. Doan, and Doina Precup. Self-supervised learning of distance functions for goal-conditioned reinforcement learning. *ArXiv*, abs/1907.02998, 2019.
- [74] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [75] Huang Xiao, Michael Herman, Joerg Wagner, Sebastian Ziesche, Jalal Etesami, and Thai Hong Linh. Wasserstein Adversarial Imitation Learning. *arXiv:1906.08113 [cs, stat]*, June 2019. URL <http://arxiv.org/abs/1906.08113>. arXiv: 1906.08113.
- [76] Yunzhi Zhang, Pieter Abbeel, and Lerrel Pinto. Automatic Curriculum Learning through Value Disagreement. *arXiv:2006.09641 [cs, stat]*, June 2020. URL <http://arxiv.org/abs/2006.09641>. arXiv: 2006.09641.
- [77] Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On Learning Intrinsic Rewards for Policy Gradient Methods. *arXiv:1804.06459 [cs, stat]*, June 2018. URL <http://arxiv.org/abs/1804.06459>. arXiv: 1804.06459.
- [78] Zeyu Zheng, Junhyuk Oh, Matteo Hessel, Zhongwen Xu, Manuel Kroiss, Hado van Hasselt, David Silver, and Satinder Singh. What Can Learned Intrinsic Rewards Capture? *arXiv:1912.05500 [cs]*, July 2020. URL <http://arxiv.org/abs/1912.05500>. arXiv: 1912.05500.

A Metrics and Quasimetrics

A metric space (M, d) is composed of a set M and a metric $d : M \times M \rightarrow \mathbb{R}^+$ that compares two points in that set. Here \mathbb{R}^+ is the set of non-negative real numbers.

Definition 2. A metric $d : M \times M \rightarrow \mathbb{R}^+$ compares two points in set M and satisfies the following axioms $\forall m_1, m_2, m_3 \in M$:

$$d(m_1, m_2) = 0 \iff m_1 = m_2 \text{ (identity of indiscernibles)}$$

$$d(m_1, m_2) = d(m_2, m_1) \text{ (symmetry)}$$

$$d(m_1, m_2) \leq d(m_1, m_3) + d(m_3, m_2) \text{ (triangle inequality)}$$

A variation on metrics that is important to this paper is *quasimetrics*.

Definition 3. A quasimetric [66] is a function that satisfies all the properties of a metric, with the exception of symmetry $d(m_1, m_2) \neq d(m_2, m_1)$.

As an example, consider an MDP where the actions and transition dynamics allow an agent to navigate from any state to any other state. Let $T(s_2 | \pi, s_1)$ be the random variable for the first time-step that state s_2 is encountered by the agent after starting in state s_1 and following policy π . The time-step metric d_T^π for this MDP can then be defined as

$$d_T^\pi(s_1, s_2) := \mathbb{E} [T(s_2 | \pi, s_1)]$$

d_T^π is a quasimetric, since the action space and transition function need not be symmetric, meaning the expected minimum time needed to go from s_1 to s_2 need not be the same as the expected minimum time needed to go from s_2 to s_1 . The diameter of an MDP [39, 43] is generally calculated by taking the maximum time-step distance between over all pairs of states in the MDP either under a random policy or a policy that travels from any state to any other state in as few steps as possible.

B Optimal Transport and Wasserstein-1 Distance

The theory of optimal transport [74, 14] considers the question of how much work must be done to transport one distribution to another optimally. More concretely, suppose we have a metric space (M, d) where M is a set and d is a metric on M . See the definitions of metrics and quasimetrics in Appendix A. For two distributions μ and ν with finite moments on the set M , the Wasserstein- p distance is denoted by:

$$W_p(\mu, \nu) := \inf_{\zeta \in \mathcal{Z}(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \zeta} [d(X, Y)^p]^{1/p} \quad (10)$$

where Z is the space of all possible couplings between μ and ν . Put another way, Z is the space of all possible distributions $\zeta \in \Delta(M \times M)$ whose marginals are μ and ν respectively. Finding this optimal coupling tells us what is the least amount of work, as measured by d , that needs to be done to convert μ to ν . This Wasserstein- p distance can then be used as a cost function (negative reward) by an RL agent to match a given target distribution [75, 19].

Finding the ideal coupling (meaning finding the optimal transport plan from one distribution to the other) which gives us an accurate distance is generally considered intractable. However, if what we need is an accurate estimate of the Wasserstein distance and not the optimal transport plan (as is the case when we mean to use this distance as part of our intrinsic reward) we can turn our attention to the dual form of this distance. The Kantorovich-Rubinstein duality [74] for the Wasserstein-1 distance on a ground metric d is of particular interest and gives us the following equality:

$$W_1(\mu, \nu) = \sup_{\text{Lip}(f) \leq 1} \mathbb{E}_y \nu [f(y)] - \mathbb{E}_x \mu [f(x)] \quad (11)$$

where the supremum is over all 1-Lipschitz functions $f : M \rightarrow \mathbb{R}$ in the metric space, and the Lipschitz constant of a function f is defined as:

$$\text{Lip}(f) := \sup \left\{ \frac{|f(y) - f(x)|}{d(x, y)} \mid x, y \in M, x \neq y \right\} \quad (12)$$

That is, the Lipschitz condition of this function f (called the Kantorovich potential function) is measured according to the metric d . Recently, Jevtić [40] has shown that this dual formulation where the constraint on the potential function is a smoothness constraint extends to quasimetric spaces as well. If defined over a quasimetric space, the Wasserstein distance also has properties of a quasimetric (specifically, the distances are not necessarily symmetric).

If the given metric space is a Euclidean space ($d(x, y) = \|x - y\|_2$), the Lipschitz bound in Equation 2 can be computed locally as a uniform bound on the gradient of f .

$$W_1(\mu, \nu) = \sup_{\phi} \left(\int \phi(y) d\nu - \int \phi(x) d\mu \right) \quad (13)$$

meaning that f is the solution to an optimization objective with the restriction that $\| \nabla f(x) \| \leq 1$ for all $x \in M$. This strong bound on the dual in Euclidean space is the one that has been used most in recent implementations of the Wasserstein generative adversarial network [3, 31] to regularize the learning of the discriminator function. Such regularization has been found to be effective for stability in other adversarial learning approaches such as adversarial imitation learning [27].

Practically, the Kantorovich potential function f can be approximated using samples from the two distributions μ and ν , regularization of the potential function to ensure smoothness, and an expressive function approximator such as a neural network. A more in depth treatment of the Kantorovich relaxation and the Kantorovich-Rubinstein duality, as well as their application in metric and Euclidean spaces using the Wasserstein-1 distance we lay out above, is provided by Peyré and Cuturi [57].

Now consider the problem of goal-conditioned reinforcement learning. Here the target distribution ν is the goal-conditioned target distribution ρ_g which is a Dirac at the given goal state. Similarly, the distribution to be transported μ is the agent's goal-conditioned state distribution ρ_π .

The Wasserstein-1 distance of an agent executing policy π to the goal s_g can be expressed in a fairly straightforward manner as:

$$W_1(\rho_\pi, \rho_g) = \sum_{s \in S} \rho_\pi(s) d(s, s_g) \quad (14)$$

The above is a simplification of Equation 1, where $p = 1$ and the joint distribution is easy to specify since the target distribution ρ_g is a Dirac distribution.

C Lipschitz constant of Potential function

For a given goal s_g and all states $s_0 \in S$, recall that function f is L -Lipschitz if it follows the Lipschitz condition as follows.

$$|f(s_g) - f(s_0)| \leq L d_T^\pi(s_0, s_g) \quad \forall s_0 \in S \quad (15)$$

Proposition 4. *If transitions from the agent policy π are guaranteed to arrive at the goal in finite time and f is L -bounded in expected transitions, i.e.,*

$$\sup_{s \in S} \mathbb{E}_{\pi, P} [|f(s^0) - f(s)|] \leq L,$$

then f is L -Lipschitz.

Proof. Since $f(s_g) - f(s_0)$ is a scalar quantity, we may write $f(s_g) - f(s_0) = \mathbb{E}_{\pi, P} [f(s_g) - f(s_0)]$. Using this fact and that $P(T(s_0) < \infty) = 1$ where $T(s_0) = T^\pi(s_g | \pi, s_0)$ for notation simplicity,

the LHS of the expression above becomes a telescopic sum

$$\begin{aligned}
jf(s_g) - f(s_0) &= \mathbb{E}_{\pi, P} [f(s_g) - f(s_0)] \\
&= \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} (f(s_{t+1}) - f(s_t)) \right] \\
&= \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} jf(s_{t+1}) - f(s_t) \right].
\end{aligned}$$

Now let us assume that for all transitions (s, a, s^θ) , $\mathbb{E}[jf(s^\theta) - f(s)] \leq L$. Then

$$\begin{aligned}
\mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} jf(s_{t+1}) - f(s_t) \right] &= \mathbb{E}_{T(s_0)} \left[\mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} jf(s_{t+1}) - f(s_t) \middle| T(s_0) \right] \right] \\
&= \mathbb{E}_{T(s_0)} \left[\sum_{t=0}^{T(s_0)-1} L \right] \\
&= L \mathbb{E}_{T(s_0)} [T(s_0)] \\
&= L d_T^\pi(s_0, s_g),
\end{aligned}$$

showing that $jf(s_g) - f(s_0) \leq L d_T^\pi(s_0, s_g)$ as desired. \square

D Proofs of Claims

The Bellman optimality condition gives us the following optimal distance to goal:

$$d_T(s, s_g) = \begin{cases} 0 & \text{if } s = s_g \\ 1 + \min_{a \in \mathcal{A}} \sum_{s^\theta \in \mathcal{S}} P(s^\theta | s, a) d_T(s^\theta, s_g) & \text{otherwise} \end{cases} \quad (16)$$

Proposition 1. A lower bound on the value of any state under a policy π can be expressed in terms of the time-step distance from that state to the goal: $V(s_0 | s_g) \geq \gamma^{d_T^\pi(s_0, s_g)}$.

Proof.

$$V^\pi(s | s_g) = \mathbb{E} \left[\gamma^{T(s_g | \pi, s)} \right] \geq \gamma^{d_T^\pi(s, s_g)} \quad \forall s \in \mathcal{S}$$

where the inequality follows as a consequence of Jensen's inequality and the convex nature of the value function. \square

Proposition 2. If the transition dynamics are deterministic, the policy that maximizes expected return is the policy that minimizes the time-step metric ($\pi^* = \pi^*$).

Proof. Consider the value of a state s given goal s_g . If the transitions are deterministic and the agent policy π is deterministic (as is the case for the optimal policy), then the time to reach the goal satisfies $\text{Var}(T(s_g | \pi, s)) = 0$, implying that Δ_{Jensen} vanishes and therefore

$$V^\pi(s | s_g) = \gamma^{d_T^\pi(s, s_g)}.$$

Since $\gamma \in [0, 1)$, V^π is monotonically decreasing with d_T^π

$$\arg \max_{\pi} V^\pi(s | s_g) = \arg \min_{\pi} d_T^\pi(s, s_g) \quad \forall s \in \mathcal{S}$$

That is, in the deterministic transition dynamics scenario, $\pi^* = \pi^*$. \square

Proposition 3. For a given policy π , the Wasserstein distance of the state visitation measure of that policy from the goal state distribution ρ_g under the ground metric d_T^π can be written as

$$W_1^\pi(\rho_\pi, \rho_g) = \mathbb{E}_{s_0} \left[h(d_T^\pi(s_0, s_g)) + \frac{\gamma}{1-\gamma} (\Delta_{Jensen}^\pi(s_0) - 1) \right] \quad (6)$$

where h is an increasing function of d_T^π .

Proof. The first step of the proof is to obtain an analytical expression for the expected distance to the goal after t steps as a function of the expected distance at $t = 0$. To reduce the notation burden, denote $T(s_0) = T(s_g | \pi, s_0)$ and let $s_t(s_0)$ be the state after t steps conditional on some starting state s_0 where actions are taken according to π . We have excluded s_g and π from the notation since they are fixed for the purpose of this proposition. Using the law of total expectation we have that for every initial s_0

$$\mathbb{E}_{s_t} [d(s_t(s_0), s_g)] = \mathbb{E}_{T(s_0)} [\mathbb{E}_{s_t} [d(s_t(s_0), s_g) | T(s_0)]] = \mathbb{E}_{T(s_0)} [\max(T(s_0) - t, 0)],$$

Now, by expanding the definition of $\rho_\pi(s | s_g)$ in equation 5, exchanging the order of summation, and using the previous equation we may write

$$\begin{aligned} W_1^\pi(\rho_\pi, \rho_g) &= \sum_{s \in \mathcal{S}} \sum_{t=0}^{\gamma} (1 - \gamma) \gamma^t \mathbb{E}_{s_0} [P(s_t = s | \pi, s_g)] d_T^\pi(s, s_g) \\ &= \mathbb{E}_{s_0} \left[(1 - \gamma) \sum_{t=0}^{\gamma} \gamma^t \mathbb{E}_{s_t} [d(s_t(s_0), s_g) | s_0] \right] \\ &= \mathbb{E}_{s_0} \left[\mathbb{E}_{T(s_0)} \left[(1 - \gamma) \sum_{t=0}^{\gamma} \gamma^t \max(T(s_0) - t, 0) \middle| s_0 \right] \right] \end{aligned}$$

Standard but tedious algebraic manipulations given in Lemma 1 in the Appendix show that

$$\sum_{t=0}^{\gamma} (1 - \gamma) \gamma^t \max(T(s_0) - t, 0) = T(s_0) - \frac{\gamma}{1-\gamma} (1 - \gamma^{T(s_0)}).$$

Combining the two identities above we arrive at

$$\begin{aligned} W_1^\pi(\rho_\pi, \rho_g) &= \mathbb{E}_{s_0} \left[\mathbb{E}_{T(s_0)} \left[T(s_0) - \frac{\gamma}{1-\gamma} (1 - \gamma^{T(s_0)}) \middle| s_0 \right] \right] \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) - \frac{\gamma}{1-\gamma} (1 - \mathbb{E}[\gamma^{T(s_0)} | s_0]) \right] \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) + \frac{\gamma}{1-\gamma} \gamma^{d(s_0, s_g)} - \frac{\gamma}{1-\gamma} (1 - \mathbb{E}[\gamma^{T(s_0)} | s_0] + \gamma^{d(s_0, s_g)}) \right] \quad (17) \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) + \frac{\gamma}{1-\gamma} \gamma^{d(s_0, s_g)} + \frac{\gamma}{1-\gamma} (\Delta_{Jensen}^\pi(s_0) - 1) \right]. \end{aligned}$$

To finalize the proof, we only need to show that the function $h(\mu) = \mu + \frac{\gamma}{1-\gamma} \gamma^\mu$ is monotonically increasing for every $\gamma \in [0, 1)$. This is a standard calculus exercise that we show in Lemma 2 in Appendix E. \square

Theorem 1. If the transition dynamics are deterministic, the policy that minimizes the Wasserstein distance over the time-step metrics in a goal-conditioned MDP (see equation 5) is the optimal policy.

Proof. Proposition 2 shows that the Jensen gap vanishes for the optimal policy of an MDP with deterministic transitions and that it minimizes the expected distance from start for all initial states. Proposition 3, on the other hand, implies that when the Jensen gap vanishes, the Wasserstein distance is monotonically increasing in the expected distance from the start. Together, the two propositions show that π^* minimizes the Wasserstein distance. \square

Algorithm 1: AIM + HER

Input: Agent policy π_θ , discriminator f_ϕ , environment env , number of Epochs N , number of time-steps per epoch K , policy update period k , discriminator update period m , episode length T , replay buffer (for HER), smaller replay buffer (for discriminator)

```
1 Initialize discriminator parameters  $\phi$ ;  
2 Initialize policy parameters  $\theta$ ;  
3 for  $n = 0, 1, \dots, N - 1$  do  
4    $t = 0$ ;  
5   goal_reached = True;  
6   while  $t < K$  do  
7     if goal_reached or episode_over then  
8       Sample goal  $s_g \sim \sigma(G)$ ;  
9       Sample start state  $s \sim \rho_0(S)$ ;  
10      goal_reached = False;  
11      episode_over = False;  
12       $t_{start} = K$ ;  
13     end  
14     Sample action  $a \sim \pi_\theta(j_s, s_g)$ ;  
15      $s^0 = env.step(a)$ ;  
16     if  $s^0 = s_g$  then  
17       | goal_reached = True;  
18     end  
19     // end episode if goal not reached in  $T$  steps  
20     if  $t - t_{start} = T$  then  
21       | episode_over = True;  
22     end  
23     Add  $(s, a, s^0, s_g, goal\_reached)$  to replay buffer and smaller replay buffer;  
24     if goal_reached or episode_over then  
25       | Add hindsight goals to both buffers;  
26     end  
27     // Update policy parameters  $\theta$  every  $k$  steps  
28     if  $t \% k = 0$  then  
29       | Sample tuples  $(s, a, s^0, s_g, goal\_reached)$  from replay buffer;  
30       | Get intrinsic reward (Equation 9);  
31       | Update policy parameters  $\theta$  using any off-policy learning algorithm;  
32     end  
33     // Update discriminator parameters  $\phi$  every  $m$  steps  
34     if  $t \% m = 0$  then  
35       | Sample tuples  $(s, a, s^0, s_g, goal\_reached)$  from smaller replay buffer;  
36       | Update discriminator parameters  $\phi$  using Equation 8;  
37     end  
38      $t = t + 1$ ;  
39   end  
40   Evaluate agent policy;  
41 end
```

E Auxiliary results for Proposition 3

Lemma 1. *Let T be a positive integer. Then*

$$\sum_{t=0}^{T-1} (1 - \gamma)\gamma^t \max(T - t, 0) = T \frac{\gamma}{1 - \gamma} (1 - \gamma^T).$$

Proof. Direct computation gives

$$\begin{aligned} (1 - \gamma) \sum_{t=0}^T \gamma^t \max(T - t, 0) &= (1 - \gamma) \sum_{t=0}^{T-1} \gamma^t (T - t) \\ &= (1 - \gamma) T \sum_{t=0}^{T-1} \gamma^t - (1 - \gamma) \sum_{t=0}^{T-1} t \gamma^t \end{aligned}$$

We will now simplify the two terms of the last expression. For the first one, have

$$(1 - \gamma) T \sum_{t=0}^{T-1} \gamma^t = (1 - \gamma) T \frac{1 - \gamma^T}{1 - \gamma} = T - T\gamma^T.$$

For the second one, the computations are a bit more involved

$$\begin{aligned} (1 - \gamma) \sum_{t=0}^{T-1} t \gamma^t &= (1 - \gamma) \gamma \sum_{t=1}^{T-1} t \gamma^{t-1} \\ &= (1 - \gamma) \sum_{t=1}^{T-1} \gamma \frac{d}{d\gamma} \gamma^t \\ &= \gamma (1 - \gamma) \frac{d}{d\gamma} \sum_{t=0}^{T-1} \gamma^t \\ &= \gamma (1 - \gamma) \frac{d}{d\gamma} \frac{1 - \gamma^T}{1 - \gamma} \\ &= \frac{\gamma}{(1 - \gamma)} \left(T \gamma^{T-1} (1 - \gamma) + (1 - \gamma^T) \right) = T \gamma^T + \frac{\gamma}{(1 - \gamma)} (1 - \gamma^T). \end{aligned}$$

When combining the two simplified expressions the terms with $T\gamma^T$ will cancel out, yielding the desired expression. \square

Lemma 2. *The function $h_\gamma(\mu) = \mu + \frac{\gamma}{1-\gamma} \gamma^\mu$ is monotonically increasing for every $\gamma \in [0, 1)$.*

Proof. We must show that $\frac{d}{d\mu} h_\gamma(\mu) > 0$ for every $\gamma \in [0, 1)$ and every $\mu > 0$. Computing the derivative directly we obtain

$$\frac{d}{d\mu} h_\gamma(\mu) = 1 + \frac{\log(\gamma) \gamma^{\mu+1}}{1 - \gamma}.$$

Thus, it will suffice to show that the second term above is greater than -1. For this purpose, first note that $\log(\gamma) \gamma^{\mu+1} > \log(\gamma)$ since $\gamma < 1$. Now, we use the fact that $\log(\gamma) < 1 - \gamma$ for $\gamma < 1$. This can be verified noting that $1 - \gamma$ is the tangent line to the concave curve $\log(\gamma)$ and the curves meet at $\gamma = 1$. And therefore $\log(\gamma)/(1 - \gamma) > -1$. Putting these observations together,

$$\frac{d}{d\mu} h_\gamma(\mu) = 1 + \frac{\log(\gamma) \gamma^{\mu+1}}{1 - \gamma} > 1 + \frac{\log(\gamma)}{1 - \gamma} > 1 - 1 = 0,$$

concluding the proof. \square

F Grid World Experiments

Basic experiment The environment is a 10 × 10 grid with 4 discrete actions that take the agent in the 4 cardinal directions unless blocked by a wall or the edge of the grid. The agent policy is learned using soft Q-learning [32], with an entropy coefficient of 0.1 and a discount factor of $\gamma = 0.99$. We do not use hindsight goals for this experiment, and use a single buffer with size 5000 for both the policy as well as the discriminator training. The results are discussed in the main text. The compute used to conduct these experiments was a personal laptop with an Intel i7 Processor and 16 GB of RAM.

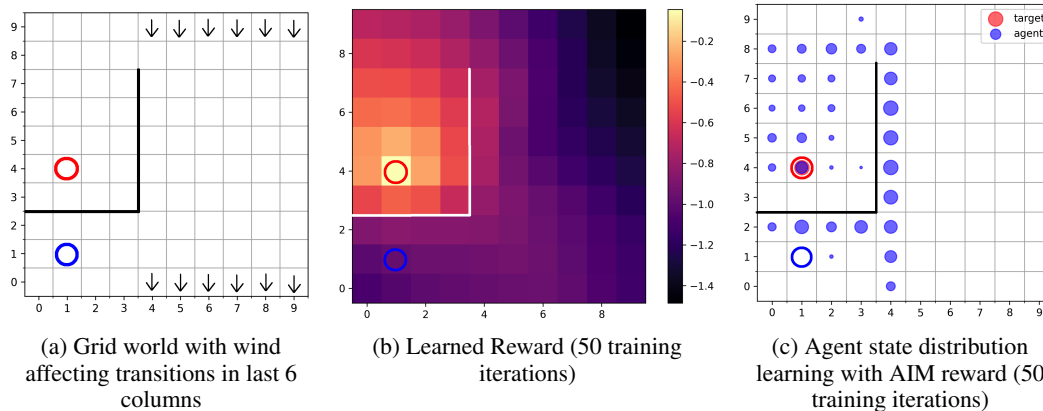


Figure 4: Windy grid world (Figure 4a) experiments. The columns with arrows at the top and bottom have stochastic and asymmetric transitions induced by wind blowing from the top. Learned reward function (Figure 4b). Reward at each state of the grid world after training for 50 iterations with AIM. Hollow red circle indicates the goal state. White lines indicate the walls the agent cannot transition through. The agent’s state visitation (Figure 4c): The hollow blue circle indicates agent’s start state. The hollow red circle is the goal. Blue bubbles indicate relative time the agent’s policy causes it to spend in respective states. Black lines indicate walls.

Additional experiments We conducted variations from the basic experiment in the grid world to show that AIM and its novel regularization can learn a reward function which guides the agent to the goal even in the presence of stochastic transitions as well as transitions where the state features vary wildly from one step to the next.

First, we evaluate AIM’s ability to learn in the presence of stochastic and asymmetric transitions in a windy version (Figure 4a) of the above grid world. Transitions in the last six columns of the grid are affected by a wind blowing from the top. Actions that try to move upwards only succeed 60% of the time, and actions attempting to move sideways cause a transition diagonally downwards 40% of the time. Movements downwards are unaffected. The rest of the experiment is carried out in the same way as above, but with 128 hidden units in the hidden layer of the agent’s Q function approximator (the reward function architecture is unchanged from the previous experiment). In Figure 4 we see that AIM learns a reward function that is still useful and interpretable, and leads to a policy that can confidently reach the goal, regardless of these stochastic and asymmetric transitions. Notice the effect of the stochastic transitions in the increased visitation in the sub-optimal states in the bottom two rows of column number 4.

The next experiment tests what happens when the transition function causes the agent to jump between states where the state features vary sharply. As an example consider a toroidal grid world, where if an agent steps off one side of the grid it is transported to the other side. The distance function here should be smooth across such transitions, but might be hampered by the sharp change in input features. In Figure 5 we see show the policy and reward for a 10 × 10 toroidal grid world with start state at (2, 2) and goal at (7, 7). Transitions are deterministic but wrap around the edges of the grid as described above: a **down** action in row 0 will transport the agent to the same column but row 9. The start and the goal state are set up so that there are multiple optimal paths to the goal. The entropy maximizing soft Q-learning algorithm should take these paths with almost equal probability. From Figure 5 it is evident that AIM learns a reward function that is smooth across the actual transitions in the environment and allows the agent to learn a Q-function that places near equal mass on multiple trajectories.

Finally, we compare learning with AIM to the baselines mentioned in Section 6. RND, SMiRL, and MC were implemented and debugged on the grid world domain with a goal that is easier to reach before being used on the Fetch robot tasks. Hyper-parameters for the algorithms in both domains were determined through sweeps. In the Fetch domains, the hyperparameters for all three new baselines were decided on through sweeps on the FetchReach task, similar to how they were evaluated for AIM and the other baselines.

Figure 6 shows the results of executing these additional baselines on the grid world domain we use to motivate AIM. All the plots are taken after the techniques have had the same number of training iterations. However none of the baselines reach the goal even after providing additional time. We show the negative L2 distance to goal as a reward in the grid world domain to highlight that the DiscoRL [50] objective should not be considered equivalent to an oracle of the distance to goal. Note that RND (Figure 6c) explores most of the larger room early on, and then converges to the state distribution seen in the figure when it does not encounter the task reward. The SMiRL reward encourages the agent to minimize surprise, and the policy trained with this reward keeps the agent in the bottom left near its start state (Figure 6d).

G Statistical Analysis of the Results on Fetch Robot Tasks

To compare the performance of each method with statistical rigor, we used a repeated measures ANOVA design for binary observation where an observation is successful if an agent reaches the goal within an episode. We then conducted a Tukey test to compare the effects of each method, i.e., the estimated odds of reaching the goal given the algorithm. The goal of the statistical analysis presented here is twofold

1. Separate the uncertainty on the performance of each method from the variation due to random seeds.
2. Adjust the probability of making a false discovery due to multiple comparisons. This extra step is necessary to avoid detecting a large fraction of falsely “significant” differences since typical tests are designed to control the error rate of only one experiment.

The data for statistical analysis comes from $N_{\text{episodes}} = 100$ evaluation episodes per each one of $N_{\text{seeds}} = 6$ seeds. For all environments but FetchReach, these data is collected after 1 million environment interactions; and for FetchReach it is taken after 2000 interactions.

The repeated measures ANOVA design is formulated as a mixed effects generalized linear model and fitted separately for each one of the four environments

$$\begin{aligned}
 y_{ijk} &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(p_{ij}), & k &\in \{1, \dots, N_{\text{episodes}}\} \\
 \text{logit}(p_{ij}) &= r_{\text{seed}_i} + \beta_{\text{algorithm}_j}, & i &\in \{1, \dots, N_{\text{seeds}}\}, j \in \{1, \dots, N_{\text{algorithms}}\} \\
 r_{\text{seed}_i} &\stackrel{\text{iid}}{\sim} \text{Normal}(0, \sigma^2)
 \end{aligned}$$

The variation due to the seed effects is measured by σ^2 , whereas the uncertainty about the odds of reaching the goal using each algorithm is measured by the standard errors of the coefficients $\beta_{\text{algorithm}_j}$.

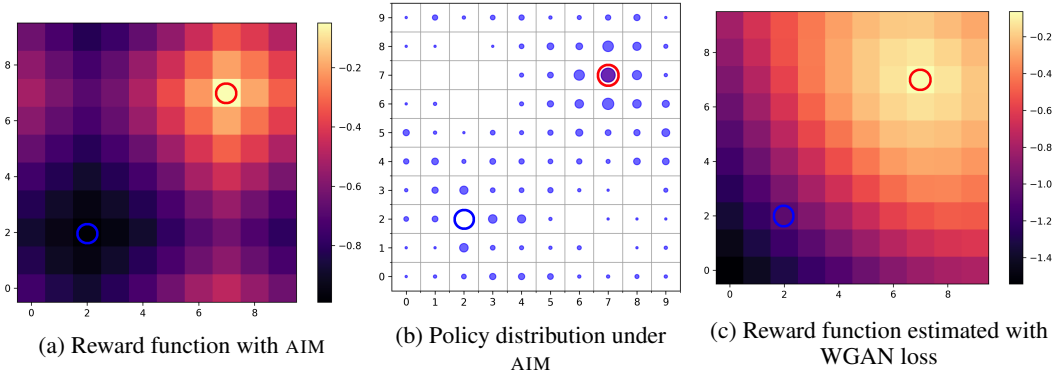


Figure 5: The reward function (Figure 5a) learned with AIM and subsequent policy distribution (Figure 5b) in a toroidal grid world, where an agent can transition from one edge of the grid across to the other. The hollow blue circle denotes the start state and the hollow red circle is the goal state. The reward function respects the sharp transitions from one end of the grid to the other. Conversely, if the reward function is learned using the WGAN objective [31] (Figure 5c), it does not respect the environment dynamics.

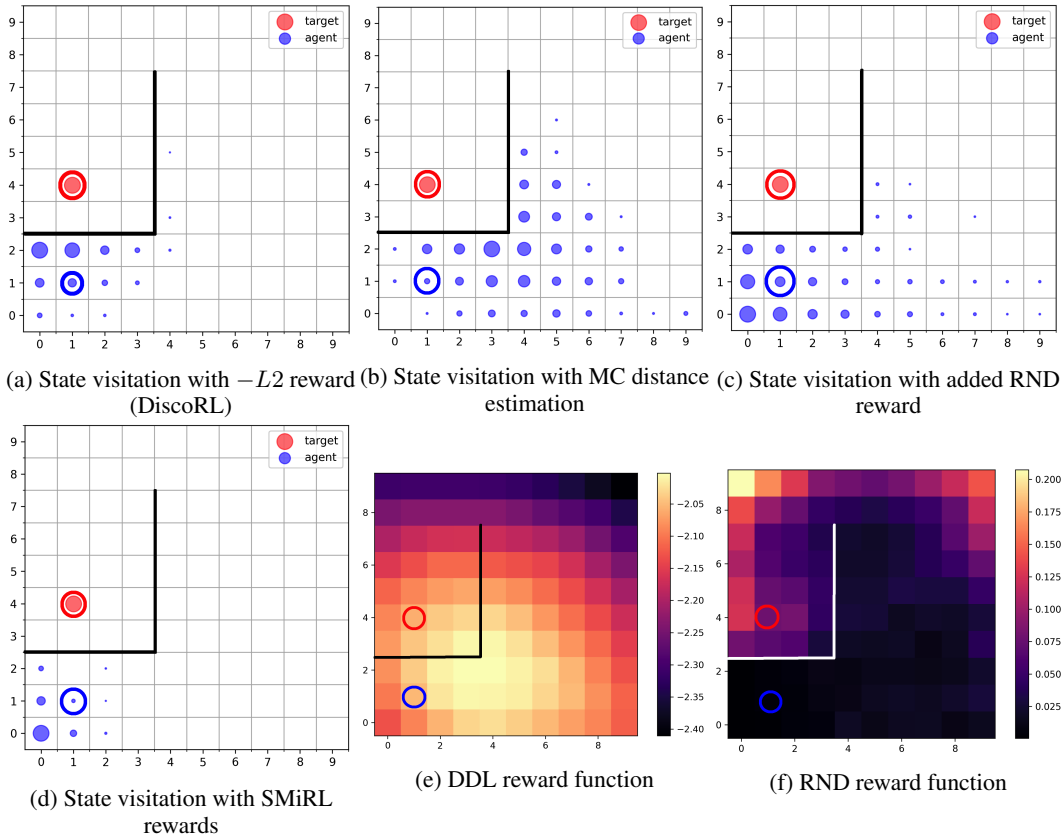


Figure 6: The state of the state visitation and reward functions for the new baselines. For comparison, Figure 2a shows the state visitation of policy trained using AIM. All algorithms are compared after 100 training iterations.

The Tukey test evaluates all null hypotheses $H_0: \beta_{\text{algorithm}_j} = \beta_{\text{algorithm}_{j^0}}$ for all combinations of j, j^0 . To adjust for multiple comparisons each Tukey tests uses the Holm method. Since we are also doing a Tukey test for each environment, we further apply a Bonferroni adjustment with a factor of four. These types of adjustments are fairly common for dealing with multiple comparison in the literature of experimental design; the interested reader may consult [49].

The results, shown in Table 1, signal strong statistical evidence of the improvements from using the AIM learned rewards. In three of the four environments AIM and AIM+ R have similar odds of reaching the goal as the dense shaped reward (H_0 is not rejected,) and in all four environments AIM and AIM+ R have statistically significant higher odds of reaching the goal than the sparse reward (H_0 is rejected and β is higher.)

Contrast	Slide	Push	PickAndPlace	Reach
$\beta_{\text{AIM+R}} - \beta_{\text{HER+dense}}$	0.34 (0.14)	-1.74 (0.77)	-0.10 (0.45)	*-3.43 (0.34)
$\beta_{\text{AIM}} - \beta_{\text{HER+dense}}$	0.21 (0.14)	-2.19 (0.75)	*-1.50 (0.37)	*-5.01 (0.35)
$\beta_{\text{AIM+R}} - \beta_{\text{HER+sparse}}$	*0.69 (0.13)	*5.32 (0.35)	*4.71 (0.33)	*4.75 (0.25)
$\beta_{\text{AIM}} - \beta_{\text{HER+sparse}}$	*0.57 (0.13)	*4.86 (0.30)	*3.31 (0.19)	*3.17 (0.24)

Table 1: Results of the Tukey test on the evaluation of Fetch tasks. The table entries are log odds ratios with standard deviations shown in parentheses. Positive values mean that AIM or AIM+R perform better than the method with negative sign in the contrast and viceversa. Asterisks mark statistical significance at 95%. If there is no asterisk, then H_0 is not rejected in which case the differences could be due to random chance.

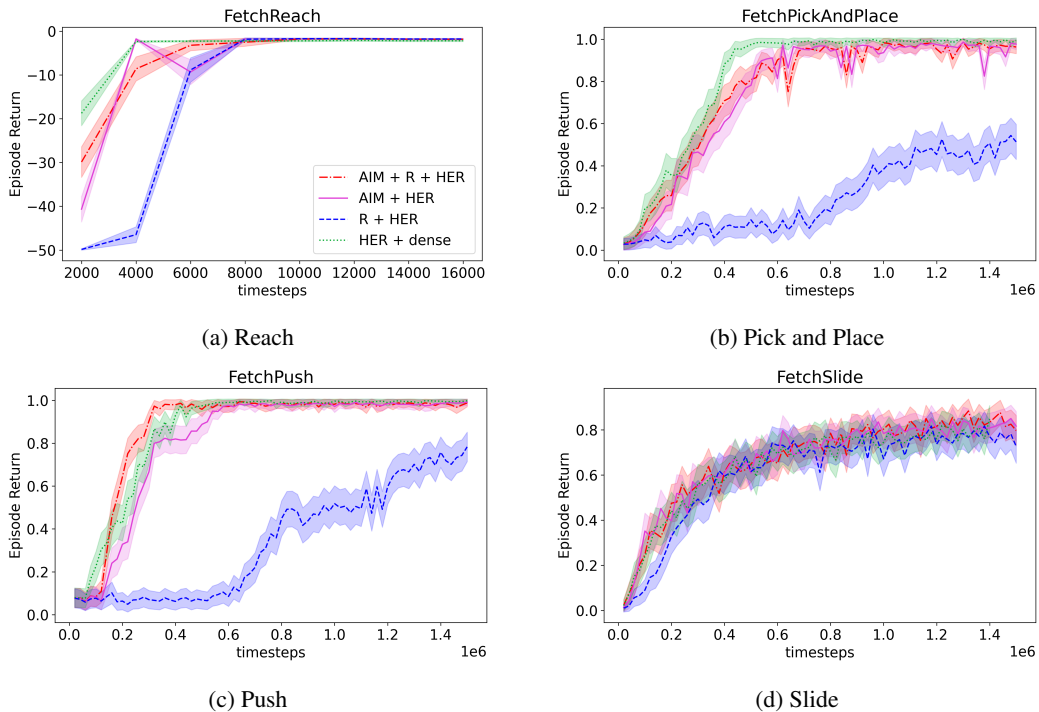


Figure 7: Comparing [AIM + HER] with an additional baseline which also uses the external task reward [AIM + R + HER]. The additional grounding provided by the external task reward allows the agent’s learning to accelerate even further.

H Details of Experiments on Fetch Robot

The Fetch robot domain in OpenAI gym has four tasks available for testing. They are named Reach, Push, Slide, and Pick And Place. The Reach task is the simplest, with the goal being the 3-d coordinates where the end effector of the robot arm must be moved to. The Push task requires pushing an object from its current position on the table to the given target position somewhere else on the table. Slide is similar to Push, except the coefficient of friction on the table is reduced (causing pushed objects to slide) and the potential targets are over a larger area, meaning that the robot needs to learn to hit objects towards the goal with the right amount of force. Finally, Pick And Place is the task where the robot actuates it’s gripper, picks up an object from its current position on the table and moves it through space to a given target position that could be at some height above the table. The goal space for the final three tasks are the required position of the object, and the goal the current state represents is the current position of that object.

Next, we note the hyperparameters used for various baselines as well as our implementation. The names of the hyperparameters are as specified in the stable baselines repository and used in the RL Zoo [59] codebase which we use for running experiments. Both the stable baselines repository and RL Zoo are available under the MIT license. These experiments were run on a compute cluster with each experiment assigned an Nvidia Titan V GPU, a single CPU and 12 GB of RAM. Each run of the TD3 baseline HER + R or HER + dense required 18 hours to execute, and each run which included AIM required 24 hours to complete execution.

TD3 [25], like its predecessor DDPG [47], suffers from the policy saturating to extremes of its parameterization. Hausknecht and Stone [35] have suggested various techniques to mitigate such saturation. We use a quadratic penalization for actions that exceed 80% of the extreme value at either end, which is sufficient to not hurt learning and prevent saturation. Assuming the policy network predicts values between -1 and 1 (as is the case when using the tanh activation function),

the regularization loss is:

$$L_a = \frac{1}{N} \sum_{i=1}^N [\max(j\pi_\theta(s_i) - 0.8, 0)]^2$$

where N is the mini-batch size and s_i is the state for the i^{th} transition in the batch.

The other modification made to the stable baselines code is to use the Huber loss instead of the squared loss for Q-learning.

For evaluation, in the Reach domain the agent policy is evaluated for 100 episodes every 2000 steps. For the other three domains, the experiment is run for 1 million timesteps, and evaluated at every 20,000 steps for 100 episodes.

H.1 TD3 and HER (R + HER)

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.95
random_exploration	0.3
target_policy_noise	0.2
learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	10
gradient_steps	10
τ (target policy update rate)	0.05

H.2 Dense reward TD3 and HER (dense + HER)

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.95
random_exploration	0.3
target_policy_noise	0.2
learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	100
gradient_steps	200
policy_delay	5
τ (target policy update rate)	0.05

H.3 TD3 and HER with AIM (AIM + HER) and (AIM + R + HER)

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.9
random_exploration	0.3
target_policy_noise	0.2
learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	100
gradient_steps	200
disc_train_freq	100
disc_steps	20
τ (target policy update rate)	0.1