
Plan Your Target and Learn Your Skills: State-Only Imitation Learning via Decoupled Policy Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 State-only imitation learning (SOIL) enables agents to learn from massive demon-
2 strations without explicit action or reward information. However, previous methods
3 attempt to learn the implicit state-to-action mapping policy directly from state-only
4 data, which results in ambiguity and inefficiency. In this paper, we overcome this
5 issue by introducing hyper-policy as sets of policies that share the same state tran-
6 sition to characterize the optimality in SOIL. Accordingly, we propose Decoupled
7 Policy Optimization (DPO) via explicitly decoupling the state-to-action mapping
8 policy as a state transition predictor and an inverse dynamics model. Intuitively,
9 we teach the agent to plan the target to go and then learn its own skills to reach.
10 Experiments on standard benchmarks and a real-world driving dataset demonstrate
11 the effectiveness of DPO and its potential of bridging the gap between reality and
12 simulations of reinforcement learning.

13 1 Introduction

14 Imitation learning offers a way to train an intelligent agent from demonstrations by mimicking the
15 expert’s behaviors without constructing hand-crafted reward functions [13, 17]. The corresponding
16 methods normally require the expert demonstrations include information of both states and actions.
17 Unfortunately, the action information is not always accessible from many real-world demonstration
18 resources, e.g., online video recordings of car driving or sports. Thus a natural desire to take advantage
19 of these massive and valuable resources motivates the study of state-only imitation learning (SOIL),
20 also known as learning from observations (LfO) [24]. Analogy to human beings, SOIL is a more
21 intuitive way to approach imitation by only matching the expert’s state sequences without having
22 explicit knowledge of the exact actions.

23 A wide range of algorithms have been proposed to solve SOIL by matching the state sequence of the
24 expert [22, 23, 25]. However, the action agnostic setting in SOIL makes it challenging to determine
25 the optimal action because of the partial observability of the expert demonstrations that multiple
26 policies could be chosen to match the same expert state sequence. Thus learning a state-to-action
27 policy is implicit, leading to a less efficient modeling of the explicit information from demonstrations,
28 and in result could cause suboptimality.

29 To this end, in this paper, we introduce the concept of *hyper-policy* denoting a *family* of policies that
30 share the same state transition. Based on that, instead of recovering the expert *policy*, we characterize
31 the optimality in SOIL by finding the expert *hyper-policy*. The proposed method is called decoupled
32 policy optimization (DPO), which separates the policy into two modules: an expert state transition
33 predictor that finds the optimal *hyper-policy*, followed by an inverse dynamics model that builds the
34 executable *policy* to deliver actions. Intuitively, the expert state transition predictor predicts the target,

35 while the inverse dynamics model enables the agent to learn its own skills to reach the target. DPO
 36 takes the advantage of such a decoupled structure by separately learning two kinds of data: (1) the
 37 expert state transition that is directly accessible in the demonstration; (2) the action to be performed
 38 which should be obtained by interacting with the environment.

39 To ensure the benefit of DPO, these two modules should work coherently to provide accurate foresight
 40 for targets and corresponding skills. To achieve this, we regularize the state transition predictor to
 41 prevent the model from predicting non-neighboring states via multi-step and cycle training style.
 42 Further, to improve the learning efficiency by encouraging the agent to reach the expert states, we
 43 augment reward and apply policy gradient to DPO with additional generative adversarial objective.

44 Experiments on standard benchmarking tasks show the advantage of the decoupled structure and
 45 the higher efficiency of DPO. We also evaluate DPO on a real-world driving dataset with state-only
 46 demonstrations, and the result shows that DPO can learn driving behaviors closer to human drivers
 47 when compared with baseline methods.

48 2 Preliminaries

49 **Markov Decision Process.** We consider a γ -discounted infinite horizon Markov decision process
 50 (MDP) as a tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \rho_0, r, \gamma \rangle$, where \mathcal{S} is the set of states, \mathcal{A} represents the action space,
 51 $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is environment dynamics distribution, $\rho_0 : \mathcal{S} \rightarrow [0, 1]$ is the distribution of
 52 the initial state s_0 , and $\gamma \in [0, 1]$ is the discount factor. The agent makes decisions through a policy
 53 $\pi(a|s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ and receives rewards $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$.

54 **Occupancy Measure.** The concept of occupancy measure (OM) [10] is proposed to characterize
 55 the statistical properties of a certain policy interacting with an MDP. Specifically, the state OM is
 56 defined as the time-discounted cumulative stationary density over the states under a given policy π :
 57 $\rho_\pi(s) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi)$. Following such a definition we can define different OM:

- 58 a) State-action OM: $\rho_\pi(s, a) = \pi(a|s)\rho_\pi(s)$
- 59 b) State transition OM: $\rho_\pi(s, s') = \int_{\mathcal{A}} \rho_\pi(s, a)\mathcal{T}(s'|s, a) da$
- 60 c) Joint OM: $\rho_\pi(s, a, s') = \rho_\pi(s, a)\mathcal{T}(s'|s, a)$

61 **Imitation Learning from State-Only Demonstrations.** Imitation learning (IL) [13] studies the
 62 task of learning from demonstrations (LfD), which aims to learn a policy from expert demonstrations
 63 without getting access to the reward signals. The expert demonstrations typically consist of expert
 64 state-action pairs. General IL objective minimizes the state-action OM discrepancy:

$$\pi^* = \arg \min_{\pi} \mathbb{E}_{s \sim \rho_{\pi_E}^*} [\ell(\pi_E(\cdot|s), \pi(\cdot|s))] \Rightarrow \arg \min_{\pi} \ell(\rho_{\pi_E}(s, a), \rho_{\pi}(s, a)) , \quad (1)$$

65 where ℓ denotes some distance metric. For example, GAIL [10] chooses to minimize the JS divergence
 66 $D_{JS}(\rho_{\pi_E}(s, a) || \rho_{\pi}(s, a))$, and AIRL [5] utilizes the KL divergence $D_{KL}(\rho_{\pi_E}(s, a) || \rho_{\pi}(s, a))$ instead,
 67 which corresponds to a maximum entropy solution with the recovered reward [17]. However, for the
 68 scenario studied in this paper, the action information is absent in state-only demonstrations, known as
 69 state-only imitation learning (SOIL) or learning from observations (LfO) problems, where the action
 70 spaces between the expert and the agent can even be different. Such challenges prevent applying
 71 typical IL solutions. An existing method for this problem is to instead optimize the discrepancy of
 72 the state transition OM with the state-to-action policy $\pi(a|s)$ [23]:

$$\pi^* = \arg \min_{\pi} [\ell(\rho_{\pi_E}(s, s'), \rho_{\pi}(s, s'))] . \quad (2)$$

73 However, the solution to this problem is ambiguous since there is no one-to-one correspondence
 74 between $\rho(s, s')$ and π as we will show in the following section. As such, the optimality of SOIL
 75 should be reconsidered.

76 3 Methodology

77 3.1 Characterizing the Optimality in SOIL

78 In standard IL tasks, when the expert actions are accessible in demonstrations, perfectly imitating the
 79 expert policy corresponds to matching the state-action OM due to the one-to-one correspondence

80 between π and $\rho_\pi(s, a)$ [10, 21]. However, such correspondence is not applicable for the state
 81 transition OM matching in SOIL.

82 **Proposition 1.** *Suppose Π is the policy space and $\mathcal{P} = \{\rho : \rho \geq 0\}$ is a feasible set of OM, then a
 83 policy $\pi \in \Pi$ corresponds to one state transition OM $\rho_\pi \in \mathcal{P}$. However, a state transition OM $\rho \in \mathcal{P}$
 84 can correspond to more than one policy in Π .*

85 The proof can be found in Appendix B. As a result, if we choose to optimize a state-to-action mapping
 86 policy, then the optimal solution to Eq. (2) is ambiguous. The ambiguity also comes from the fact
 87 that Eq. (2) does not correspond to a maximum policy entropy solution as in normal IL tasks (see
 88 Appendix C for details). Therefore, a state-to-action mapping function may be too implicit for
 89 matching the state sequence, which could cause training instability and lead to sub-optimal policies.
 90 In that case, we must find a one-to-one corresponding solution to solve SOIL explicitly and efficiently.
 91 Before continuing, we introduce the definition of *hyper-policy*.

92 **Definition 1.** *A hyper-policy Ω is a set of policies such that for any $\pi_1, \pi_2 \in \Omega$, we have $\rho_{\pi_1}(s, s') =$
 93 $\rho_{\pi_2}(s, s')$.*

94 Then by definition, there is a one-to-one correspondence between the hyper-policy Ω and the state
 95 transition OM $\rho_\Omega(s, s')$. Similar to the normal state-to-action mapping policy, a hyper-policy Ω can
 96 be regarded as a state-to-state mapping function $h_\Omega(s'|s)$ which predicts the state transition such that
 97 for any $\pi \in \Omega$:

$$h_\Omega(s'|s) = \frac{\rho_\Omega(s, s')}{\int_{\tilde{s}} \rho_\Omega(s, \tilde{s}) d\tilde{s}} = \int_a \pi(a|s) \mathcal{T}(s'|s, a) da. \quad (3)$$

98
 99 **Proposition 2.** *Suppose the state transition predictor is defined as in Eq. (3) and Γ is its space,
 100 $\mathcal{P} = \{\rho : \rho \geq 0\}$, then a hyper-policy state transition predictor $h_\Omega \in \Gamma$ corresponds to one state
 101 transition OM $\rho_\Omega \in \mathcal{P}$; and a state transition OM $\rho \in \mathcal{P}$ only corresponds to one hyper-policy state
 102 transition predictor such that $h_\rho = \rho(s, s') / \int_{\tilde{s}} \rho(s, \tilde{s}) d\tilde{s}$.*

103 Therefore, we find a one-to-one correspondence between the optimization term $\rho(s, s')$ and a practical
 104 target $h_\Omega(s'|s)$, which indicates that we do not have to infer the expert actions under state-only
 105 demonstrations but only need to recover the state transition predictor of the hyper-policy Ω_E :

$$\arg \min_{\Omega} [\ell(\rho_{\Omega_E}(s, s'), \rho_\Omega(s, s'))] \Rightarrow \arg \min_{h_\Omega} \mathbb{E}_{s \sim \Omega} [\ell(h_{\Omega_E}(s'|s), h_\Omega(s'|s))]. \quad (4)$$

106 However, SOIL still requires to learn a policy to interact with the MDP environment to match the
 107 state transition OM of the expert. This is achievable since we do not have to recover the expert policy
 108 π_E exactly but can learn any policy $\pi \in \Omega_E$ according to Eq. (4).

109 3.2 Policy Decoupling

110 To construct an unambiguous objective for SOIL, we define hyper-policy and solve the problem by
 111 finding the state transition predictor of the expert hyper-policy. Intuitively, this tells the agent the
 112 *target* that the expert will reach without informing any feasible *skill* that require the agent to learn
 113 itself. Therefore, to recover a $\pi \in \Omega_E$, we can construct an inverse dynamics such that

$$\pi = \underbrace{\mathcal{T}_\pi^{-1}}_{\text{Inverse dynamics}} \left(\underbrace{\mathcal{T}(\pi_E)}_{\text{Expert state transition predictor}} \right). \quad (5)$$

114 Formally, the expert policy can be decoupled as

$$\begin{aligned} \pi_E(a|s) &= \int_{s'} \mathcal{T}(s'|s, a) \pi_E(a|s) ds' = \int_{s'} \frac{\rho_{\pi_E}(s, a, s')}{\rho_{\pi_E}(s)} ds' = \int_{s'} \frac{\rho_{\pi_E}(s, s') I_{\pi_E}(a|s, s')}{\rho_{\pi_E}(s)} ds' \\ &= \int_{s'} h_{\pi_E}(s'|s) I_{\pi_E}(a|s, s') ds'. \end{aligned} \quad (6)$$

115

116 Notice that both the state transition predictor h and the inverse dynamics model I is policy dependent.
 117 Nevertheless, recall that the optimality in SOIL only requires us to recover $\pi \in \Omega_E$, we do not have
 118 to learn about I_{π_E} but just one feasible skill $I(a|s, s')$. Then a policy can be recovered by

$$\pi = \mathbb{E}_{\underbrace{s' \sim h_{\Omega_E}(s'|s)}_{\text{target}}} \left[\underbrace{I(a|s, s')}_{\text{skill}} \right]. \quad (7)$$

119 Here the inverse dynamics model I offers an
 120 arbitrary *skill* to reach the expected *target*
 121 state provided by the state transition predictor h . In
 122 fact, it does not depend on the hyper-policy
 123 Ω_E but a sampling policy π_B to construct $I =$
 124 I_{π_B} . We only need a mild requirement for π_B
 125 that it covers the support of $\rho_{\Omega_E}(s, s')$ so that
 126 the learned I can provide a possible action to
 127 achieve the target state. In both experiments and
 128 theoretical analysis we show that this require-
 129 ment alleviates the dependence on the inverse
 130 dynamics. Furthermore, if the environment and
 131 the expert policy are both deterministic (which
 132 is usually the case in real-world scenarios such
 133 as robotics), the state transition is a single-point
 134 distribution (or known as the Dirac delta func-
 135 tion), and we can simply model h as a determinis-
 136 tic function. By decoupling the policy, which is
 137 a state-to-action mapping function, as a state-to-
 138 state mapping function (the transition predictor)

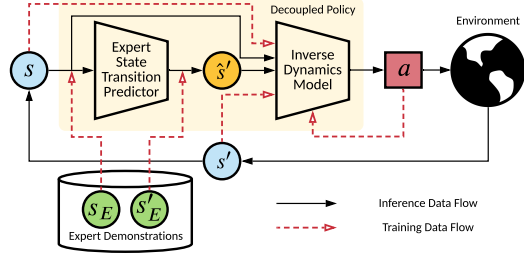


Figure 1: The architecture of Decoupled Policy Optimization (DPO), which consists of an expert state transition predictor (to plan where to go) followed by an inverse dynamics model (to decide how to reach).

139 **State Transition Predictor.** In practice, we construct a parameterized expert state transition predictor h_ψ which predicts the subsequent state of the expert taking the input as a current state $\hat{s}' = h_\psi(s)$.
 140 The state transition predictor models the explicit information of the expert, and it can be learned from the demonstration data only. Thence, we implement Eq. (4) as a KL divergence minimization:

$$\min_{\psi} \mathbb{E}_{(s, s') \sim \Omega_E} [\text{D}_{\text{KL}}(h_{\Omega_E}(s'|s) \| h_\psi(s'|s))], \quad (8)$$

143 which can be optimized in a supervised manner. Specifically, we sample state transitions (s, s') from the expert demonstrations \mathcal{D} and optimize the L2 loss:

$$\mathcal{L}_\psi^h = \mathbb{E}_{(s, s') \sim \mathcal{D}} [\|s' - h_\psi(s)\|^2]. \quad (9)$$

145 **Inverse Dynamics Model.** Knowing where to go is not enough since the agent has to interact with the environment to reach the target. This can be achieved via an inverse dynamics model, which predicts the action given two consecutive states. Formally, let the ϕ -parameterized inverse dynamics model I_ϕ take input the state pair and predict the feasible action to achieve the state transition: $\hat{a} = I_\phi(s, s')$. Intuitively, we want the inverse dynamics to learn from possible transitions sampled by the agent. Recall that we only need the support of learned $I(a|s, s')$ of the sampling policy covers the support of the expert state transition OM, from which we can infer at least one possible action. Hence, we can optimize the KL divergence between the inverse dynamics of a sampling policy π_B and I_ϕ :

$$\min_{\phi} \mathbb{E}_{(s, s') \sim \pi_B} [\text{D}_{\text{KL}}(I_{\pi_B}(a|s, s') \| I_\phi(a|s, s'))], \quad (10)$$

154 and we can choose to optimize L2 loss in a supervised manner by sampling from the replay buffer \mathcal{B} :

$$\mathcal{L}_\phi^I = \mathbb{E}_{(s, a, s') \sim \mathcal{B}} [\|a - I_\phi(s, s')\|^2]. \quad (11)$$

155 In our implementation, both the state predictor and the inverse dynamics can be constructed as Gaussian distributions similar to a normal stochastic policy, thus encouraging exploration.

157 3.3 Tackling Compounding Error Challenges

158 In our formulation, we have decoupled the state-to-action mapping policy as a state-to-state mapping function and a state-pair-to-action mapping function. Unfortunately, the compounding error problem exists such that the agent cannot reach where it plans due to the fitting errors of these two parts.

161 **Theorem 1** (Error Bound of DPO). *Consider a deterministic environment whose dynamics transition function $\mathcal{T}(s, a)$ is deterministic and L -Lipschitz. Assume the ground-truth state transition $h_{\Omega_E}(s)$ is deterministic, and for each policy $\pi \in \Pi$, its inverse dynamics I_π is also deterministic and*

164 *C-Lipschitz. Then for any state s , the distance between the desired state s'_E and reaching state s'*
 165 *sampled by the decoupled policy is bounded by*

$$\|s' - s'_E\| \leq LC \|h_{\Omega_E}(s) - h_{\psi}(s)\| + L \|I_{\pi_B}(s, s') - I_{\phi}(s, \hat{s}')\|, \quad (12)$$

166 *where π_B is a sampling policy that covers the state transition support of the expert hyper-policy and*
 167 *$\hat{s}' = h_{\psi}(s)$ is the predicted next state.*

168 The proof can be found in Appendix B, where we also
 169 induce a similar error bound for rollout with a state-
 170 to-action policy as BCO [22] to show the advantage
 171 of the decoupled structure. From Theorem 1 we
 172 know that the compounding error can be enlarged
 173 due to each part’s fitting error, where the first term
 174 corresponds to the error of predicted states and the
 175 second term indicates whether the agent can reach
 176 where it plans to. To alleviate the error, we further propose regularization on these two modules.

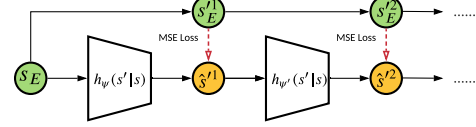


Figure 2: Multi-step optimization. Given an expert state s_E , h_{ψ} predicts the next possible state \hat{s}'^1 , which is further fed to a target network $h_{\psi'}$ to predict the following sequence. The total loss computes the MSE loss along the state sequence.

177 3.3.1 Regularization on Target Planning

178 One major problem is that the state transition predictor may suggest non-neighboring states instead
 179 of predicting one-step reachable states. To overcome this, we draw inspiration from Asadi et al. [2]
 180 and Edwards et al. [3], and regularize state transition predictor to prevent the model from predicting
 181 non-neighboring states via multi-step and cycle training style.

182 **Multi-Step Optimization.** We first explain the details of the multi-step optimization objective.
 183 This idea is motivated by Asadi et al. [2], which optimizes a multi-step outcome by executing a
 184 sequence of actions in the dynamics model. Here we optimize the state sequence instead. As shown
 185 in Fig. 2, given an expert state s_E , h_{ψ} predicts the next possible state \hat{s}' that the expert will reach; the
 186 predicted state is then fed into the predictor to output the predicted two-step state \hat{s}'' . As such, the
 187 multi-step training loss is the L2 loss computed along the k -step outcome sequence:

$$\mathcal{L}_{\psi}^{h,ms} = \mathbb{E}_{(s, \{s_E^{i'}\}_{i=1}^k) \sim \mathcal{D}} \left[\|s_E^{1'} - h_{\psi}(s)\|^2 + \sum_{i=2}^k \|s_E^{i'} - h_{\psi'}(s^{i-1})\|^2 \right]. \quad (13)$$

188 Intuitively, such a regularization makes the state prediction
 189 \hat{s}' close to the expert state distribution in order to make
 190 accurate long step predictions. It is worth noting that the
 191 gradient of the cascading state transition predictors should
 192 be dropped since we already have an accurate input at each
 193 time step, and for each training step, we only update the
 194 first one. We use a target network $h_{\psi'}$ in practice.

195 **Cycle Training Style.** Another way to regularize the
 196 transition predictor’s output to a neighboring state is to
 197 keep an additional function to ensure the cycle consistency,
 198 which is also an important technique in [3]. In particular,
 199 as illustrated in Fig. 3, given an expert state s_E , we take
 200 the predicted state \hat{s}' and s_E into the inverse dynamics
 201 get the action a , then we train an additional forward dynamics
 202 that takes the input (s_E, a) and gets a forward next state \tilde{s}' :

$$\begin{aligned} \mathcal{L}_{\omega}^M &= \mathbb{E}_{(s, a, s') \sim \mathcal{B}} [\|s' - M_{\omega}(s, a)\|^2] \\ \mathcal{L}_{\psi}^{h,cycle} &= \mathbb{E}_{(s, s') \sim \mathcal{D}} [\|s' - h_{\psi}(s)\|^2 + \|h(s) - M_{\omega}(s, I_{\phi}(s, s'))\|^2]. \end{aligned} \quad (14)$$

203 In other words, the cycle training scheme provides a regularization on h_{ψ} to make predictions
 204 consistent with the forward dynamics model.

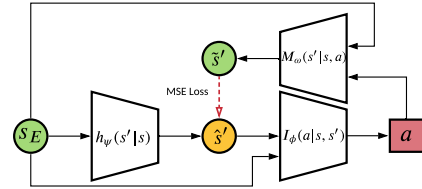


Figure 3: Cycle training style. Given an expert state s_E , $I_{\phi}(s, s')$ takes input the predicted state \hat{s}' and s_E to get the execution action a , then an additional forward dynamics model M_{ω} is used to simulate one step rollout using (s_E, a) and get a forward next state \tilde{s}' . The total loss computes the MSE loss between the two predicted states.

205 **3.3.2 Efficient Skills Learning via Decoupled Policy Gradient**

206 In previous sections, we have mentioned that learning to reach a specific place requires the data-
 207 collecting policy to cover the support of the expert hyper-policy. This is easy to achieve on simple
 208 low-dimensional tasks, but may not be satisfied in high-dimensional continuous environments. To
 209 this end, we encourage the agent to approach those state transitions from the expert’s hyper-policy
 210 Ω_E by minimizing the JS divergence of the state transition occupancy using a state-to-action mapping
 211 policy $D_{JS}(\rho_{\pi_E}(s, s'), \rho_{\pi}(s, s'))$. This can be done by producing informative rewards via GAN-like
 212 methods [10, 23], and updating the decoupled policy with policy gradients (PG).

213 In detail, we construct a parameterized discriminator $D_{\omega}(s, s')$ to compute the reward $r(s, a) \triangleq$
 214 $r(s, s')$ as $\log D_{\omega}(s, s')$ and the decoupled policy served as the generator. In addition, since we
 215 decouple the policy as two parameterized modules, i.e., a state transition predictor and an inverse
 216 dynamics model, then by chain rule, the PG for the decoupled policy can be accomplished by

$$\begin{aligned} \nabla \mathcal{L}_{\phi, \psi}^{\pi} &= \mathbb{E}_{\pi} [Q(s, a) \nabla_{\phi, \psi} \log \pi_{\phi, \psi}(a|s)] \\ &= \mathbb{E}_{\pi} \left[Q(s, a) \int_{s'} (\nabla_{\psi} \log h_{\psi}(s'|s) + \nabla_{\phi} \log I_{\phi}(a|s, s')) ds' \right], \end{aligned} \quad (15)$$

217 where Q is the state-action value function; the first term is the gradient for updating the state transition
 218 predictor; and the second term is for the inverse dynamics model. Thus, the optimization for both
 219 the state transition predictor and the inverse dynamics model can augment the supervised learning
 220 objectives with any PG-based learning algorithms (e.g., TRPO, PPO, SAC). As the training proceeds,
 221 the agent will sample more transition data around Ω_E , and thus the support of the sampling policy
 222 will progressively cover the support of $\rho_{\Omega_E}(s, s')$.

223 **3.4 Overall Algorithm**

224 By combining the idea of generative adversarial training, we obtain our final algorithm, composed
 225 with three essential parts: the state transition predictor h used for predicting the possible future
 226 states sampled by the expert; the inverse dynamics model I used for inferring the possible actions
 227 conditioned on two adjacent states; and the discriminator D used for offering intermediate reward
 228 signals for training the decoupled policy $\pi = I(h)$. The overall objective of DPO is

$$\mathcal{L}_{\phi, \psi}^{\pi, h, I} = \lambda_G \mathcal{L}_{\phi, \psi}^{\pi} + \lambda_h \mathcal{L}_{\psi}^h + \lambda_I \mathcal{L}_{\phi}^I, \quad (16)$$

229 where λ_G, λ_h and λ_I are hyperparameters for trading off the training among each loss. In practice,
 230 we try less than ten combinations for these parameters as shown in Appendix D.3, and we directly
 231 optimize $\mathcal{L}_{\phi, \psi}^{\pi}$ instead of iterative training. The detailed algorithm is summarized in Appendix A.
 232 Besides, it is worth noting that both the inverse dynamics model and the state transition predictor
 233 can be pre-trained, where we optimize \mathcal{L}_{ψ}^h using the state-only demonstration and optimize \mathcal{L}_{ϕ}^I using
 234 samples collected by a randomized agent.

235 **4 Related Work**

236 State-only imitation learning (SOIL) endows the
 237 agent with the ability to learn from expert states. Al-
 238 though lacking the expert decision information, most
 239 of the previous works still optimize a state-to-action
 240 mapping policy to match the expert state transition
 241 distribution. For example, Torabi et al. [22] used a model-based approach to apply behavioral
 242 cloning to state-only demonstrations, while Torabi et al. [23] employed a similar structure to GAIL
 243 to match the state transition distribution. Yang et al. [25] analyzed the optimization gap between
 244 SOIL and naive IL and introduced a mutual information term to narrow it. Huang et al. [11] applied
 245 SOIL on autonomous driving tasks by decoupling the policy into a neural decision module and a
 246 non-differentiable execution module in a hierarchical way.

247 Our work decouples the state-to-action policy into two modules. However, both the inverse dynamics
 248 model and the state transition predictor have been widely used by many previous works on RL and IL
 249 tasks. For instance, Torabi et al. [22] and Guo et al. [6] trained an inverse dynamics model to label the

Table 1: Comparison between different methods.

| Method | Inverse Dynamics | State Predictor | Decoupled Policy | Task |
|---------------|------------------|-----------------|------------------|------|
| BCO [22] | ✓ | ✗ | ✗ | SOIL |
| GAIPO [23] | ✗ | ✗ | ✗ | SOIL |
| IDDM [25] | ✗ | ✗ | ✗ | SOIL |
| OPOLO [27] | ✓ | ✗ | ✗ | SOIL |
| PID-GAIL [11] | ✗ | ✗ | ✓ | IL |
| QSS [3] | ✓ | ✓ | ✓ | RL |
| SAIL [16] | ✓ | ✓ | ✗ | IL |
| DPO (Ours) | ✓ | ✓ | ✓ | SOIL |

Table 2: Eventual performance against different methods on 6 easy-to-hard continuous control benchmarks. The means and the standard deviations are evaluated over more than 5 random seeds.

| | InvertedPendulum | InvertedDoublePendulum | Hopper | Walker2d | HalfCheetah | Ant |
|--------------|-----------------------|--------------------------|------------------------|-------------------------|--------------------------|-----------------------|
| Random | 25.28 ± 5.53 | 78.28 ± 10.73 | 13.09 ± 0.10 | 7.07 ± 0.13 | 74.48 ± 12.39 | 713.59 ± 203.92 |
| BCO | 1000.00 ± 0.00 | 415.04 ± 148.46 | 1430.16 ± 398.81 | 261.36 ± 25.17 | -13.66 ± 149.94 | 397.79 ± 239.16 |
| GAIfo | 1000.00 ± 0.00 | 7818.07 ± 1778.67 | 3068.10 ± 26.32 | 3865.20 ± 341.90 | 8953.35 ± 1079.41 | 5122.29 ± 807.19 |
| GAIfo-DP | 1000.00 ± 0.00 | 7305.01 ± 1591.23 | 3031.84 ± 152.13 | 4003.06 ± 241.34 | 8675.42 ± 807.29 | 5535.9 ± 62.74 |
| DPO (w/o PG) | 1000.00 ± 0.00 | 3545.70 ± 738.16 | 629.84 ± 344.07 | 334.23 ± 85.42 | -472.00 ± 132.81 | -196.96 ± 124.26 |
| DPO (w PG) | 1000.00 ± 0.00 | 7846.40 ± 1541.20 | 3165.72 ± 68.44 | 4407.53 ± 266.72 | 10501.96 ± 438.01 | 5338.48 ± 107.2 |
| Expert (SAC) | 1000.00 ± 0.00 | 9358.87 ± 0.10 | 3402.94 ± 446.48 | 5639.32 ± 29.97 | 13711.64 ± 111.47 | 5404.55 ± 1520.49 |

state-only demonstrations with inferred actions. Nair et al. [19] proposed a method for manipulating ropes to match a single human-specified image sequence, in which an inverse dynamics model is trained in a self-supervised manner and used to generate control signals. Kimura et al. [14] utilized a state transition predictor to fit the state transition probability in the expert data, which is further used to compute a predefined reward function. Liu et al. [16] constructed a policy prior using the inverse dynamics and the state transition predictor, but the policy prior was only used for regularizing the policy network. However, as shown in this paper, the policy can be exactly decoupled as these two parts, which can be uniformly optimized through policy gradient without keeping an extra policy. Edwards et al. [3] estimated $Q(s, s')$ for RL tasks which employs a similar policy form as Eq. (6) and updates the state transition predictor through a deterministic policy gradient similar to DDPG [15]. To sort out the difference between these methods and ours, we summarize the key factors in Tab. 1.

5 Experiments

We conduct four sets of experiments to investigate the following research questions:

RQ1 Is decoupled learning structure superior than state-to-action structure on SOIL tasks?

RQ2 Does DPO achieve higher efficiency or better performance than baselines on SOIL tasks?

RQ3 Can an agent reach where it plans and does the proposed regularization help mitigate the compounding error?

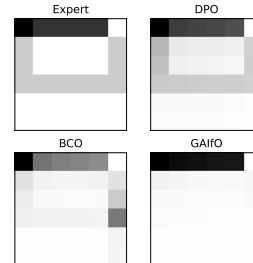
RQ4 How can DPO be applied on real-world data?

To answer RQ1, we conduct toy experiments with a simple 2D grid world environment and compare both qualitative and quantitative results of DPO against BCO and GAIfo. Regarding RQ2, we empirically evaluate DPO on easy-to-hard continuous control benchmarking tasks. And for RQ3, we evaluate the difference between the predicted states that the agent plans to reach and the consecutive state that the agent actually reaches in the environment for the proposed regularization. Finally, we try to imitate real-world traffic surveillance recordings in a simulated environment to investigate RQ4, which shows the potential of using real-world data for human behavior simulation. Due to the space limit, we leave experiment details, additional results and ablation study in Appendix.

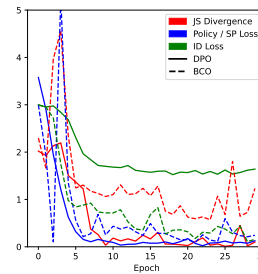
5.1 Understanding the Decoupled Structure

In this paper we design decoupled policy optimization (DPO) to perform SOIL tasks, and in previous sections we propose that the key technical contribution of DPO is the decoupled structure of policy that models the explicit state transition information and the latent action information from demonstrations, which solves the ambiguity and enhances the learning efficiency. Therefore, in this set of experiments, we aim to demonstrate how DPO is superior than state-to-action policy methods (RQ1). We first generate expert demonstrations in a 2D 6×6 grid world environment, in which the agent starts at the upper left corner and aims to reach the upper right corner. In each grid the agent has $k \times 4$ actions, which means that the agent has k possible actions to reach the neighboring block and in our experiment we choose $k = 5$ to enlarge the action space.

The density of the expert trajectories and the trajectories sampled by different methods are shown in Fig. 4(a). We show that both BCO and GAIfo have troubles in directly learning the implicit action from state-only behaviors. Notably, GAIfo only imitates the major trajectory and omit the other choice and BCO also stuck in the middle right. By contrast, DPO recovers the expert demonstrations



(a) Rollout density.



(b) Loss curves.

Figure 4: Toy example.

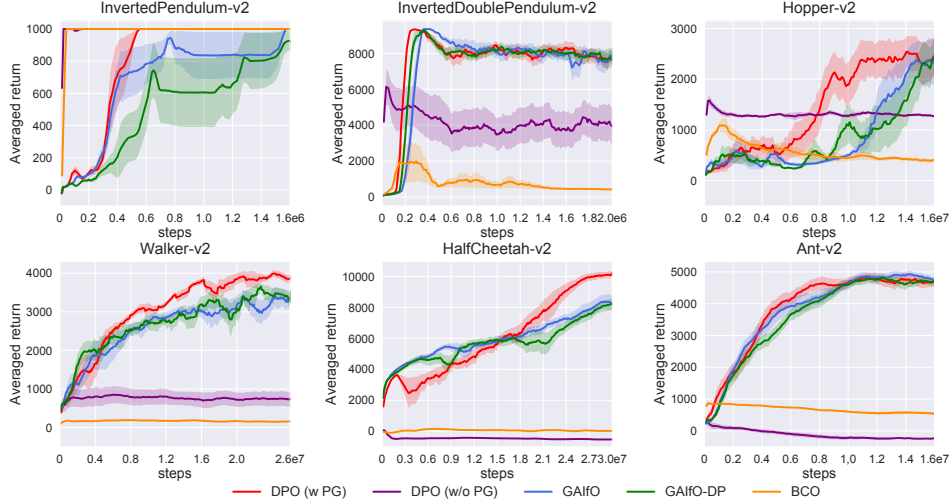


Figure 5: Learning curves on 6 easy-to-hard continuous control benchmarks, where the solid line and the shade represent the mean and the standard deviation of the averaged return over more than 5 random seeds. We pre-train BCO and DPO for 50k steps and show it in figures.

295 much better, benefiting from the decoupled structure that first determining the target and then taking
 296 the action to achieve it. To further illustrate the learning efficiency advantage of DPO, we illustrate the
 297 JS divergence curves of DPO and BCO during training in Fig. 4(b). Besides, we show the policy loss
 298 for BCO, the state predictor (SP) loss for DPO, and the inverse dynamics (ID) loss for both methods.
 299 Except that the JS divergence of DPO decreases more quickly than BCO, it is also observable that
 300 DPO relies less on the inverse dynamics than BCO, since the inverse dynamics loss of DPO converges
 301 to a higher level. We further provide a theoretic analysis of the dependence on inverse dynamics with
 302 BCO and DPO in Appendix B.

303 5.2 Comparative Evaluations

304 We compare the qualitative results of DPO against other baseline methods on easy-to-hard continuous
 305 control benchmarking environments (RQ2), including InvertedPendulum, InvertedDoublePendulum,
 306 Hopper, Walker2d, HalfCheetah and Ant. In each environment, besides GAIfo and BCO, we also
 307 evaluate GAIfo with decoupled policy (denoted as GAIfo-DP). For DPO we compare the reward
 308 augmented version of DPO (denoted as DPO w PG)¹ with the supervised learning version of DPO,
 309 i.e., $\lambda_G = 0$ (denoted as DPO w/o PG). For fairness, we re-implement all the algorithms based on a
 310 Pytorch code framework² and adopt Soft Actor-Critic (SAC) [7] as the RL learning algorithm for
 311 GAIfo and DPO. For all environments, we first train an SAC agent to collect 4 state-only expert
 312 trajectories and then train agents with such data. All algorithms are evaluated by a deterministic
 313 policy. The eventual results are summarized in Tab. 2, and the learning curves are shown in Fig. 5.
 314 It is worth noting that for DPO, we choose the best performance among the experiments that use
 315 multi-step or cycle regularization, and we put the full experiment results in Appendix D.

316 One can easily observe that on simple environments, BCO is able to achieve a good performance,
 317 and GAIfo also does well on harder tasks. Even so, DPO can still gain the best or comparable
 318 performance against its counterparts. Particularly, without augmented reward, DPO is able to reach
 319 the optimality with the highest sample efficiency on simple tasks like InvertedPendulum. By contrast,
 320 on higher-dimensional tasks such as Hopper, Walker2d, HalfCheetah and Ant, it is difficult to
 321 construct accurate inverse dynamics that covers the support of the expert hyper-policy from scratch.
 322 However, by combining generative adversarial policy gradients, the agent finally recovers a good
 323 policy from the expert hyper-policy. This is particularly evident on HalfCheetah where DPO behaves
 324 poorly at the beginning but improves fast as the training proceeds. Besides, as illustrated in Fig. 5,
 325 DPO benefits from better sample efficiency in most of the environments, but the improvements are
 326 limited on the hardest tasks. We think that this may be due to larger state spaces (111 dimensions
 327 for Ant) that makes it difficult to recover a good state predictor or an inverse dynamics model. In all

¹Without ambiguity we simply denote DPO for this version of algorithm in the following sections.

²https://github.com/KamyarGh/rl_swiss

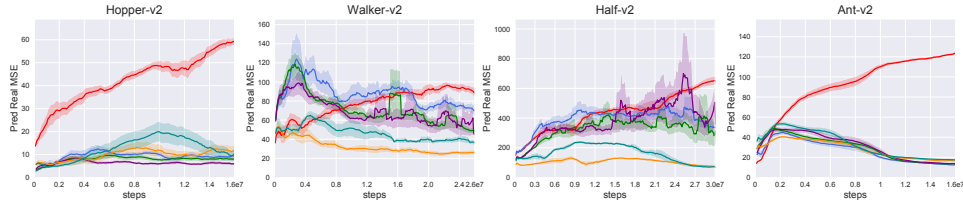


Figure 6: Compounding error of the predicted consecutive states and the real states the agent reaches when rollout in the environments.

328 experiments, GAIfO-DP achieves similar results as GAIfO, indicating that the network structure does
 329 not count much for the performance.

330 5.3 Compounding Error Reduction

331 In this section, we aim to study whether the agent can effectively reach the target as it plans (RQ3).
 332 Therefore, we analyze the distance of the reaching states and the predicted consecutive states, and
 333 draw the mean square error (MSE) along the training procedure in Fig. 6. We also compare our
 334 regularization including multi-step optimization (denoted as M.S.- k , where k is the number of rollout
 335 steps) and cycle training style (denoted as Cycle). Note that DPO needs at least 1-step rollout for
 336 training the state transition predictor. As shown in Fig. 6, the agent still has gaps to get to where it
 337 plans to, and the mismatch always deteriorates on harder tasks. Combining regularization can always
 338 achieve lower compounding error, and the cycle training is effective in most of the environments. In
 339 Appendix D.6, we further illustrate the correlation between the final performance and the distance.

340 5.4 Learn to Drive from Real-World Traffic Data

341 The rapid development of autonomous driving has brought a lot of demand for simulating and
 342 training an RL agent in the simulator, which requires realistic interactions with various social vehicles
 343 [26]. However, driver’s detailed actions are not easily to obtain yet we adopt SOIL from a traffic
 344 surveillance recording dataset (NGSIM I-80 [8]) that contains kinds of recorded human driving
 345 trajectories. We wish to further examine the potential of DPO for decreasing the gap between the
 346 real world and simulation (RQ4). We utilize the simulator provided by Henaff et al. [9] as our
 347 simulation platform and learn to imitate real-world driving behaviors. We compare DPO against
 348 GAIfO and BCO, and choose *Success Rate*, *Mean Distance* and *KL Divergence* as evaluation metrics.
 349 Specifically, *Success Rate* is the percentage of driving across the entire area without crashing into
 350 other vehicles or driving off the road, *Mean Distance* is the distance traveled before the episode ends,
 351 and *KL Divergence* measures the position distribution distance between the expert and the agent.

352 As shown in Tab. 3, DPO outperforms baseline methods in
 353 all three metrics while possessing higher stability. The de-
 354 coupled policy allows the state predictor to focus on match-
 355 ing the distribution of expert trajectories, thus achieving
 356 smaller deviations from the expert position distribution.
 357 Furthermore, since the policy gradient can be computed
 358 with non-differentiable inverse dynamics, we can generate
 359 stable action sequences [12, 11] by replacing the inverse
 360 dynamics model with classical controllers, which can be generalized to realistic applications.

Table 3: Performance on NGSIM I-80 driving task over 5 random seeds.

| Method | Success Rate (%) | Mean Distance (m) | KL Divergence |
|--------|----------------------------------|-----------------------------------|---------------------------------|
| BCO | 27.4 ± 1.1 | 129.8 ± 2.0 | 24.4 ± 2.2 |
| GAIfO | 77.5 ± 0.8 | 188.3 ± 1.1 | 11.5 ± 3.9 |
| DPO | 80.3 ± 0.5 | 192.7 ± 0.6 | 9.5 ± 1.8 |
| Expert | 100 | 210.0 | 0 |

361 6 Conclusion

362 In this paper, we characterize the optimality and investigate the ambiguity problem in state-only
 363 imitation learning, and accordingly propose Decoupled Policy Optimization (DPO), which splits
 364 the state-to-action mapping policy into a state-to-state mapping state transition predictor and a
 365 state-pair-to-action mapping inverse dynamics model. Furthermore, we employ regularization and
 366 generative adversarial methods to mitigate the compounding error caused by the decoupled modules.
 367 The flexibility of the decoupled architecture allows a wide range of interesting future works, such
 368 as replacing the inverse dynamics with a classic control module to produce stable control signals,
 369 learning specific skills with shared state transition and multi-task target learning with shared pre-
 370 trained skills.

371 **References**

- 372 [1] Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning.
373 In *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004)*,
374 2004.
- 375 [2] Kavosh Asadi, Dipendra Misra, Seungchan Kim, and Michel L Littman. Combating the
376 compounding-error problem with a multi-step model. *arXiv preprint arXiv:1905.13320*, 2019.
- 377 [3] Ashley D. Edwards, Himanshu Sahni, Rosanne Liu, Jane Hung, Ankit Jain, Rui Wang, Adrien
378 Ecoffet, Thomas Miconi, Charles Isbell, and Jason Yosinski. Estimating $q(s,s')$ with deep
379 deterministic dynamics gradients. In *Proceedings of the 37th International Conference on*
380 *Machine Learning, ICML 2020*, 2020.
- 381 [4] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal
382 control via policy optimization. In *International Conference on Machine Learning*, pages 49–58,
383 2016.
- 384 [5] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse
385 reinforcement learning. In *6th International Conference on Learning Representations, ICLR*
386 *2018*, 2018.
- 387 [6] Xiaoxiao Guo, Shiyu Chang, Mo Yu, Gerald Tesauro, and Murray Campbell. Hybrid reinforce-
388 ment learning with expert state sequences. In *The Thirty-Third AAAI Conference on Artificial*
389 *Intelligence, AAAI 2019*, pages 3739–3746, 2019.
- 390 [7] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
391 maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the*
392 *35th International Conference on Machine Learning, ICML 2018*, pages 1856–1865, 2018.
- 393 [8] John Halkias and James Colyar. Next generation simulation fact sheet. *US Department of*
394 *Transportation: Federal Highway Administration*, 2006.
- 395 [9] Mikael Henaff, Alfredo Canziani, and Yann LeCun. Model-predictive policy learning with
396 uncertainty regularization for driving in dense traffic. In *7th International Conference on*
397 *Learning Representations, ICLR 2019*, 2019.
- 398 [10] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in*
399 *Neural Information Processing Systems 29*, pages 4565–4573, 2016.
- 400 [11] Junning Huang, Sirui Xie, Jiankai Sun, Qirui Ma, Chunxiao Liu, Dahua Lin, and Bolei Zhou.
401 Learning a decision module by imitating driver’s control behaviors. In *Proceedings of the*
402 *Conference on Robot Learning (CoRL) 2020*, 2020.
- 403 [12] Sandy H. Huang, Nicolas Papernot, Ian J. Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial
404 attacks on neural network policies. In *5th International Conference on Learning Representations,*
405 *ICLR 2017*, 2017.
- 406 [13] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning:
407 A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- 408 [14] Daiki Kimura, Subhajit Chaudhury, Ryuki Tachibana, and Sakyasingha Dasgupta. Internal
409 model from observations for reward shaping. *arXiv preprint arXiv:1806.01267*, 2018.
- 410 [15] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval
411 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning.
412 In *4th International Conference on Learning Representations, ICLR 2016*,, 2016.
- 413 [16] Fangchen Liu, Zhan Ling, Tongzhou Mu, and Hao Su. State alignment-based imitation learning.
414 In *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- 415 [17] Minghuan Liu, Tairan He, Minkai Xu, and Weinan Zhang. Energy-based imitation learning. In
416 *20th International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2021*,
417 2021.

- 418 [18] Leland McInnes and John Healy. UMAP: uniform manifold approximation and projection for
419 dimension reduction. *CoRR*, abs/1802.03426, 2018.
- 420 [19] Ashvin Nair, Dian Chen, Pulkit Agrawal, Phillip Isola, Pieter Abbeel, Jitendra Malik, and Sergey
421 Levine. Combining self-supervised learning and imitation for vision-based rope manipulation.
422 In *2017 IEEE International Conference on Robotics and Automation, ICRA 2017*, pages 2146–
423 2153, 2017.
- 424 [20] Murray Rosenblatt. Remarks on some nonparametric estimates of a density function. *The*
425 *Annals of Mathematical Statistics*, pages 832–837, 1956.
- 426 [21] Umar Syed, Michael H. Bowling, and Robert E. Schapire. Apprenticeship learning using linear
427 programming. In *Machine Learning, Proceedings of the Twenty-Fifth International Conference*
428 *(ICML 2008)*, pages 1032–1039, 2008.
- 429 [22] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In
430 *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence,*
431 *IJCAI 2018*, pages 4950–4957, 2018.
- 432 [23] Faraz Torabi, Garrett Warnell, and Peter Stone. Adversarial imitation learning from state-only
433 demonstrations. In *Proceedings of the 18th International Conference on Autonomous Agents*
434 *and MultiAgent Systems, AAMAS '19*, pages 2229–2231, 2019.
- 435 [24] Faraz Torabi, Garrett Warnell, and Peter Stone. Recent advances in imitation learning from
436 observation. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial*
437 *Intelligence, IJCAI 2019*, pages 6325–6331, 2019.
- 438 [25] Chao Yang, Xiaojian Ma, Wenbing Huang, Fuchun Sun, Huaping Liu, Junzhou Huang, and
439 Chuang Gan. Imitation learning from observations by minimizing inverse dynamics disagree-
440 ment. In *Advances in Neural Information Processing Systems 32*, pages 239–249, 2019.
- 441 [26] Ming Zhou, Jun Luo, Julian Vilella, Yaodong Yang, David Rusu, Jiayu Miao, Weinan Zhang,
442 et al. Smarts: Scalable multi-agent reinforcement learning training school for autonomous
443 driving. In *Conference on Robot Learning*, 2020.
- 444 [27] Zhuangdi Zhu, Kaixiang Lin, Bo Dai, and Jiayu Zhou. Off-policy imitation learning from
445 observations. In *Advances in Neural Information Processing Systems 33*, 2020.

446 **Checklist**

- 447 1. For all authors...
- 448 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
449 contributions and scope? [Yes]
- 450 (b) Did you describe the limitations of your work? [Yes] See Section 3.3.
- 451 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
452 Section 6.
- 453 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
454 them? [Yes]
- 455 2. If you are including theoretical results...
- 456 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Theorem 1.
457 (b) Did you include complete proofs of all theoretical results? [Yes] See Section B.
- 458 3. If you ran experiments...
- 459 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
460 mental results (either in the supplemental material or as a URL)? [Yes]
- 461 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
462 were chosen)? [Yes] See Appendix D.3.
- 463 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
464 ments multiple times)? [Yes] We ran our results with more than 5 random seeds as said
465 in Section 5.2.
- 466 (d) Did you include the total amount of compute and the type of resources used (e.g., type
467 of GPUs, internal cluster, or cloud provider)? [No]
- 468 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 469 (a) If your work uses existing assets, did you cite the creators? [Yes] We re-implement all
470 algorithms based on an existed code base, as said in Section 5.2.
- 471 (b) Did you mention the license of the assets? [No]
- 472 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
473 We will public our codes after publication.
- 474 (d) Did you discuss whether and how consent was obtained from people whose data you’re
475 using/curating? [N/A]
- 476 (e) Did you discuss whether the data you are using/curating contains personally identifiable
477 information or offensive content? [N/A]
- 478 5. If you used crowdsourcing or conducted research with human subjects...
- 479 (a) Did you include the full text of instructions given to participants and screenshots, if
480 applicable? [N/A]
- 481 (b) Did you describe any potential participant risks, with links to Institutional Review
482 Board (IRB) approvals, if applicable? [N/A]
- 483 (c) Did you include the estimated hourly wage paid to participants and the total amount
484 spent on participant compensation? [N/A]