

Cross Domain Robot Imitation with Invariant Representation

Zhao-Heng Yin¹, Lingfeng Sun³, Hengbo Ma³, Masayoshi Tomizuka³, Wu-Jun Li²

Abstract—Animals are able to imitate each others’ behavior, despite their difference in biomechanics. In contrast, imitating other similar robots is a much more challenging task in robotics. This problem is called cross domain imitation learning (CDIL). In this paper, we consider CDIL on a class of similar robots. We tackle this problem by introducing an imitation learning algorithm based on invariant representation. We propose to learn invariant state and action representations, which align the behavior of multiple robots so that CDIL becomes possible. Compared with previous invariant representation learning methods for similar purposes, our method does not require human-labeled pairwise data for training. Instead, we use cycle-consistency and domain confusion to align the representation and increase its robustness. We test the algorithm on multiple robots in the simulator and show that unseen new robot instances can be trained with existing expert demonstrations successfully. Qualitative results also demonstrate that the proposed method is able to learn similar representations for different robots with similar behaviors, which is essential for successful CDIL.

I. INTRODUCTION

Animals are able to imitate each others’ behavior by watching their demonstrations despite their difference in biomechanics such as body length, shape, and strength [28]. However, previous robotics research suggests that imitating a similar reference robot of different embodiment and dynamics is a challenging problem, which is also termed as Cross Domain (robot) Imitation Learning (CDIL) in the literature [9]. Solving this problem can be both appealing and significant for robotics. Domain discrepancies between the expert and the agent usually exist in the practice of imitation learning, and direct imitation by Behavior Cloning (BC) [1] or Inverse Reinforcement Learning (IRL) [13] without noticing such discrepancies will harm learning performance [9]. CDIL will also make imitation learning more flexible and convenient, as this enables us to use the existing massive amount of demonstrations online to train similar robots.

Previous CDIL research usually focuses on cross domain imitation learning between a pair of robots [11], [22]. But the demonstrations can actually come from multiple sources (experts) in practice. Therefore in this paper, we consider a more general CDIL problem among a class of similar robots.

¹Zhao-Heng Yin is with the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR. zhaohengyin@gmail.com

²Wu-Jun Li is with the Department of Computer Science and Technology, Nanjing University, Nanjing 210012, PRC. liwujun@nju.edu.cn

³Lingfeng Sun, Hengbo Ma and Masayoshi Tomizuka are with the Department of Mechanical Engineering, University of California, Berkeley, Berkeley, CA 94720, USA. {lingfengsun, hengboma, tomizuka}@berkeley.edu

Code: https://github.com/zhaohengyin/irgail_example

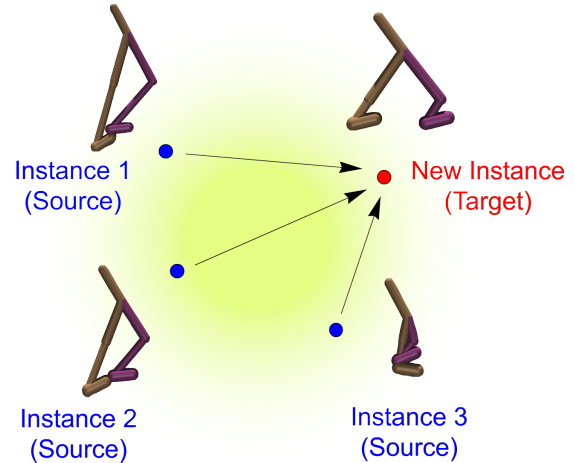


Fig. 1. Illustration of the CDIL problem studied in the paper. Each dot indicates a robot instance. Our goal is to train a new robot instance (target) with the existing expert (source) demonstrations of similar robots.

One example is shown in Figure 1, which is taken from the MuJoCo Walker problem [25]. We are provided with demonstrations of several walkers of different embodiments (leg length), and our goal is to use these demonstrations to train a new Walker instance sampled from a class of Walker robots.

The availability of demonstration from multiple experts inspires us to mine the invariant representation for behavior patterns, which can be used for CDIL. Taking Walker as an example, one can notice that all these human-like walker instances walk by alternately moving one leg forward, despite their difference in leg length. Such moving pattern is invariant and shared across the instances. This observation is formalized as the *domain-invariant subspace assumption* in recent work [30]:

Assumption The state space and the action space of similar domains can be disentangled into several independent subspaces (factors). Some of these subspaces are domain-invariant and shared by source and target tasks. Hence, we can share the task knowledge on these domain-invariant subspaces.

Our solution directly follows the assumption. We first learn a domain-invariant state and action representation space for robots and then perform imitation learning inside this invariant representation space. For example, one possible domain-invariant state representation for the Walkers is shown in Figure 2. We expect that such invariant state representation only encodes their walking behavior, and ignore the irrelevant factors. In other words, this invariant representation can be considered as a behavior prototype for these robots.

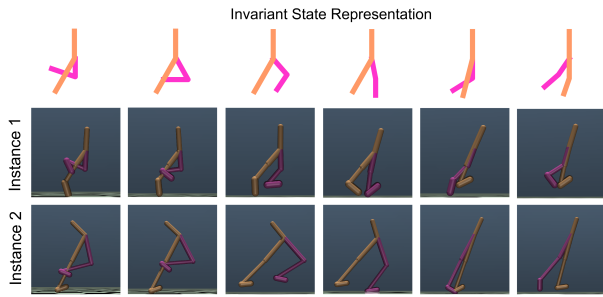


Fig. 2. Illustration of a possible domain-invariant state representation for Walkers. Despite their physical difference, the behavior of these Walkers can be extracted and described by invariant representation (above row). We propose a method to learn such invariant representation and use it to solve the CDIL problem in this paper.

There exist some representation learning methods for similar purposes [5], [10]. However, these methods require human-specified paired states from the source and target domain. This can be quite inconvenient when the number of experts increases. Moreover, specifying such pairwise states can also be difficult since it is hard to define which states can be paired in some cases [9], [29]. Our proposed method overcomes this problem by leveraging the cycle consistency structure and using domain confusion to align the representation and make it robust.

Our contributions can be summarized as follow:

- We study the CDIL problem on a class of similar robots and propose a new CDIL algorithm based on invariant representation. The proposed algorithm is more convenient since it does not depend on human-labeled pairwise data between source and target task.
- We test the proposed algorithm on various robots in the MuJoCo simulator. Experimental results show that our algorithm can outperform previous methods and generalize to unseen new robots. We also visualize to show that the proposed method is able to learn representations of similar behavior patterns of different robots.

II. RELATED WORK

A. Imitation Learning

In order to enable the robots to acquire desired behavior, one basic approach is by Reinforcement Learning (RL) [23]. However, one drawback of applying RL is the difficulty of reward design. To solve this problem, researchers propose behavior acquisition by Imitation Learning (IL) [8]. One straightforward approach of IL is Behavioral Cloning [1]. Behavioral Cloning solves the IL problem via a supervised process, in which it learns to predict expert-like action given the state. However, Behavioral Cloning algorithms may suffer from the covariate shift problem [16]. Another line of work is Inverse Reinforcement Learning (IRL) [13]. IRL algorithms propose to infer the reward function from the expert demonstration. Among IRL algorithms, one recent branch is the Adversarial Imitation Learning (AIL) [7], [26],

which trains the agent to match the expert's behavior via an adversarial process. Compared with Behavioral Cloning, AIL can succeed in various challenging control tasks [7]. In this work, our imitation learning process follows the AIL framework.

B. Cross Domain Imitation Learning

One crucial assumption of IL algorithms in the previous part is that the expert (source) and the agent (target) should be in the same domain. However, the expert demonstration may come from other similar but different domains in practice. Utilizing numerous existing demonstrations to solve the learning problem on an unseen new domain is quite appealing, and is called CDIL or Adaptive IL (ADIL) [9], [10]. One common ritual of CDIL is to map the states and actions in the source domain to functionally similar states in the agent's domain, which is called domain adaptation (translation) [29]. Such mechanism is also observed in biology [12]. Some recent research in robotics directly defines cross domain mapping by human prior knowledge [14] to train a quadruped robot, but the approach is application-specific and can be confined to particular scenarios. Besides, researchers also consider learning such mapping. Some methods learn the mapping by using pairwise data in the source and target domain [10]. Some other method utilizes auxiliary tasks and dynamics information to align the states and actions without such pairwise data [9]. However, collecting pairwise data or defining tasks on multiple domains can be tricky and difficult. Therefore, these methods can be quite inconvenient in practice. Some vision-based works also propose unsupervised domain translation [20], but these methods can misalign states in some cases [29].

We address this problem by learning invariant representations, rather than direct domain translation. Our idea is partially inspired by research in robot transfer learning [30], whose purpose is quite similar to CDIL. Some previous methods like [5] propose to learn an invariant feature space for transfer learning. However, it requires supervised pairwise data for alignment. Some other methods tackle the CDIL problem by domain confusion [22], which can also be considered as instances of invariant representation learning. We will discuss these methods in the next section. Learning invariant representations is also related to domain randomization [24], which uses various domains to make representation used by policy invariant and robust.

One limitation of our method is that it does not consider cross morphology cases (such as different number of joints). But we believe that this problem can be solved by padding state and action space in the future. Cross morphology CDIL is recently studied in some recent research [6].

C. Domain Confusion

Some CDIL methods reuse the demonstrations from different robots by getting rid of the domain-specific information in the demonstration via representation learning, and this process is usually called domain confusion [3], [4], [27].

The process is implemented by minimizing the mutual information between representation and domain label [3], [22]. To estimate the mutual information, a mutual information estimator based on the neural network called MINE is often used [2]. Besides MINE-based implementation, [15] uses a variational information bottleneck to regularize the mutual information, which is also used in recent domain adaptation method [21]. This variational information bottleneck can provide an upper bound estimation of mutual information. Some other methods like [4], [17], [29] also propose to train a domain discriminator for domain confusion. However, we find that merely using domain confusion for CDIL is not optimal in some cases, since it does not use the dynamics information to align the representation. Our method further takes this into consideration.

III. BACKGROUND AND PRELIMINARIES

A. Notations

Since similar robots usually differ from each other only on some particular physical configuration (parameters) like link length and body mass, in this paper we assume that we can use these physical configurations c to describe a robot, which is denoted as R_c . Then, a class of similar robots can be denoted as $RC = \{R_c | c \in \mathcal{C}\}$, where \mathcal{C} is the configuration space (i.e. set of possible physical configurations). For example, we can use $\mathcal{C}_{pend} = [1.5, 2]$ to characterize a class of pendulums whose lengths are between 1.5m and 2m. $R_{1.6}$ can represent a pendulum whose length is 1.6m. A similar environment characterization is also used by [11].

We formalize the control problem of a robot R_c by the Markov Decision Process (MDP) in RL, which is defined as a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$. \mathcal{S} denotes its state space. \mathcal{A} denotes its action space. \mathcal{P} denotes the transition dynamics from time step t to $t + 1$. $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function. At time step t , the robot agent observes $s_t \in \mathcal{S}$ and output an action $a_t \in \mathcal{A}$. Then, the environment evolves to s_{t+1} following the dynamics defined by \mathcal{P} . Moreover, for the transition (s_t, a_t, s_{t+1}) , we also use c_t to denote the corresponding robot's configuration. In this paper, we assume that the dimension of the robots' state space and action space are the same (the controllable joints are the same across the robots). But note that one state (action) can be of different meanings to different robots. For example, let the state space of a pendulum be the tilting angle (for simplicity). Then pendulums of different lengths should react differently to the same observed tilting angle. This is due to the difference in dynamics of different robots.

B. Generative Adversarial Imitation Learning

We use Generative Adversarial Imitation Learning (GAIL) [7] as our imitation learning framework. GAIL enables the agent to obtain expert-like behavior by encouraging it to match its behavior with the expert by fooling a discriminator network D . Such discriminator D is trained to discern expert π_E 's state-action pair from agent

π 's state-action pair, which is equivalent to minimizing the following loss:

$$-\mathbb{E}_{\pi_E} [\log(D(s_t, a_t))] - \mathbb{E}_{\pi} [\log(1 - D(s_t, a_t))].$$

Then, we can use such discriminator D to define the imitation reward for the agent. The reward for taking a_t at state s_t is defined as $r_t = -\log(1 - D(s_t, a_t))$.

C. Mutual Information Neural Estimation

Our algorithm requires mutual information estimation to learn robust invariant representation. To evaluate the mutual information, we apply the estimation method proposed by MINE [2]. This method follows from the Donsker Varadhan representation theorem, which states that given two random variables X and Z on the sample space Ω , the mutual information between X and Z is given by

$$I(X, Z) = \sup_{T: \Omega \rightarrow \mathbb{R}} \mathbb{E}_{x \sim P_{XZ}} T(x) - \log \mathbb{E}_{x \sim P_X P_Z} e^{T(x)}.$$

In the above equation, P_{XZ} denotes the joint probability distribution of X and Z , $P_X P_Z$ denotes the product of their marginal distributions. The supremum is taken over functions such that the expectation terms are finite. MINE proposes to parameterize the mapping T with a neural network T_ϕ , so

$$I_T(X, Z) = \mathbb{E}_{x \sim P_{XZ}} T_\phi(x) - \log \mathbb{E}_{x \sim P_X P_Z} e^{T_\phi(x)}$$

provides a lower bound estimation of $I(X, Z)$. We can get an approximation of $I(X, Z)$ by maximizing I_{T_ϕ} and calculate its value.

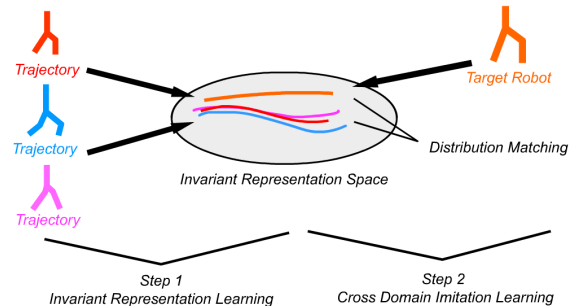


Fig. 3. Overall process of the proposed algorithm. The first step is to train an invariant representation module. The second step is cross domain imitation learning on the target robot. We train the target robot to match expert's behavior in the invariant representation space.

IV. ALGORITHM

A. Overview

We first provide the overview of our proposed learning process in Figure 3. The whole process is composed of two steps. The first step is to train an invariant representation module. This module maps the states and actions of each robot instance to the invariant representation space where similar states and actions in different robots are aligned. In the second step, we carry out imitation learning. We transform the state-action pairs of the target robot into the invariant representation space and use GAIL to train the target robot to match them with the expert demonstration.

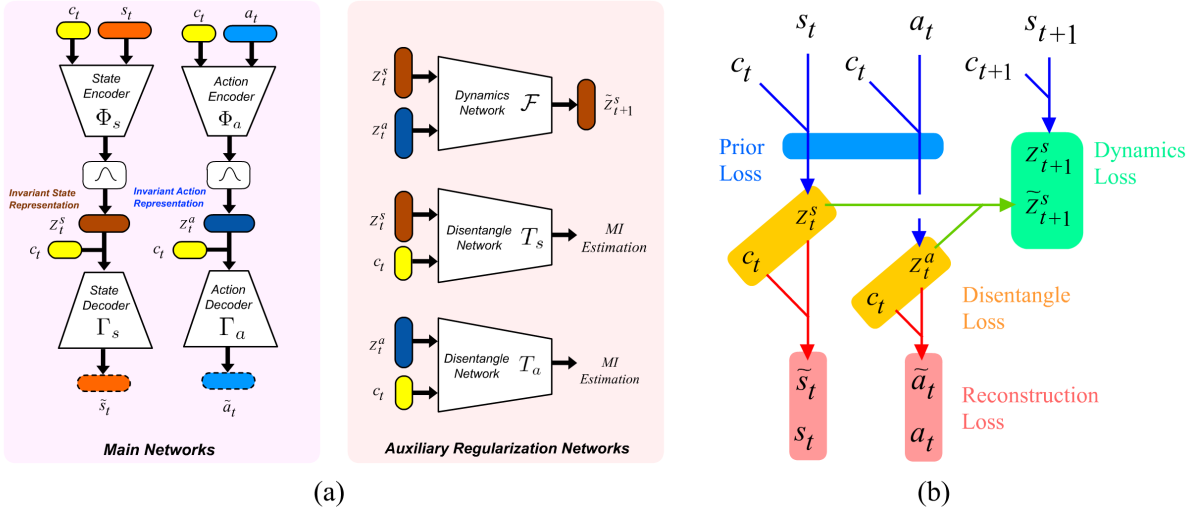


Fig. 4. (a) Structure of the proposed invariant representation module. (b) Illustration of the loss functions.

B. Invariant Representation Module

We illustrate the invariant representation module in Figure 4. It has two parts: the main part consists of encoding-decoding networks, which produce invariant representation. The other part consists of some auxiliary regularization networks to further improve the representation.

1) *Main Network*: We use two variational encoders to encode states and actions respectively. Concretely, we define the state encoder as $\Phi_s(\cdot|s_t, c_t)$, which takes current state s_t and robot configuration c_t as input, and produces a Gaussian distribution on the invariant state space \mathcal{Z}_S . The encoded invariant state representation z_t^s is then sampled from $\Phi_s(\cdot|s_t, c_t)$. The action encoder $\Phi_a(\cdot|a_t, c_t)$ is defined on the invariant action space \mathcal{Z}_A in the same way. We sample invariant action representation $z_t^a \sim \Phi_a(\cdot|a_t, c_t)$. Then, we use two decoders to ensure that the representation contains enough information about the input. The state decoder $\Gamma_s(z_t^s, c_t)$ is required to reconstruct robot's state s_t conditioned on z_t^s with its configuration c_t . The action decoder $\Gamma_a(z_t^a, c_t)$ is defined similarly. The reconstructed state and action are denoted as \tilde{s}_t and \tilde{a}_t . Then the reconstruction loss functions for state and action representation are defined as

$$\mathcal{L}_{sr} = \mathbb{E}_{\Phi_s(z|s_t, c_t)}[\|s_t - \tilde{s}_t\|^2] \quad (1)$$

and

$$\mathcal{L}_{ar} = \mathbb{E}_{\Phi_a(z|a_t, c_t)}[\|a_t - \tilde{a}_t\|^2] \quad (2)$$

respectively. To enforce Gaussian distribution, a prior loss is also used. It is defined as

$$\mathcal{L}_{kl} = \mathcal{L}_{skl} + \mathcal{L}_{akl}, \quad (3)$$

where

$$\mathcal{L}_{skl} = D_{KL}(\Phi_s(z|s_t, c_t)||p(z)), \quad (4)$$

$$\mathcal{L}_{akl} = D_{KL}(\Phi_a(z|a_t, c_t)||p(z)). \quad (5)$$

Here, D_{KL} is the Kullback-Leibler divergence. $p(z)$ is the prior Gaussian distribution $\mathcal{N}(0, 1)$.

2) *Dynamics Loss*: As is introduced before, similar robot instances share their motion (dynamics) pattern in the invariant representation space. Therefore, we also train a dynamics model in the representation space for all the robot instances to discover shared pattern. The dynamics model is a neural network $\mathcal{F} : \mathcal{Z}_S \times \mathcal{Z}_A \rightarrow \mathcal{Z}_S$, whose goal is to predict the representation of the future state given the current state and action. This prediction is denoted as $\tilde{z}_{t+1}^s = \mathcal{F}(z_t^s, z_t^a)$. The loss function is defined as

$$\mathcal{L}_{dyn} = \mathbb{E}\|\tilde{z}_{t+1}^s - z_{t+1}^s\|^2 = \mathbb{E}\|\mathcal{F}(z_t^s, z_t^a) - z_{t+1}^s\|^2. \quad (6)$$

This loss function is a cycle consistency constraint between the original state space and representation space. [29] finds that such constraint can push similar states (actions) towards each other, which is also essential for our problem.

3) *Disentangle Loss*: A necessary condition to make representation invariant is that the representation of different robots should be indistinguishable, otherwise the representation must carry robot-specific information. This condition is equivalent to minimizing the mutual information between encoded representation and robot's configuration. In our setting, both $I(z_t^s, c_t)$ and $I(z_t^a, c_t)$ should be minimized. Hence we train two MINE networks T_s and T_a to provide estimation for $I(z_t^s, c_t)$ and $I(z_t^a, c_t)$, which are written as $I_{T_s}(z_t^s, c_t)$ and $I_{T_a}(z_t^a, c_t)$ respectively. Then we define the disentangle loss for encoders as

$$\mathcal{L}_{disent} = I_{T_s}(z_t^s, c_t) + I_{T_a}(z_t^a, c_t). \quad (7)$$

The loss function for updating T_s and T_a is defined as

$$\mathcal{L}_T = -(I_{T_s}(z_t^s, c_t) + I_{T_a}(z_t^a, c_t)). \quad (8)$$

Our final training loss is defined as a weighted combination of the above loss functions:

$$\mathcal{L} = \mathcal{L}_{sr} + \mathcal{L}_{ar} + \lambda_1 \mathcal{L}_{disent} + \lambda_2 \mathcal{L}_{dyn} + \lambda_3 \mathcal{L}_{kl}. \quad (9)$$

Here, λ_1 , λ_2 and λ_3 are weighting hyperparameters. To obtain the training data, we sample transition data (s_t, a_t, s_{t+1})

on various robot instances in the given robot family using a random policy. However, applying a random policy for sampling may leave some important transitions uncovered in some tasks. Therefore, we also add rollout transition data of experts into the training dataset if such a situation happens. This is enough as expert rollouts can cover all the important transitions. Since the encoders Φ_s and Φ_a are evolving during the training process, we update T_a and T_s accordingly to track the change of the mutual information. In the implementation, we update network T_s and T_a by minimizing \mathcal{L}_T after each update step of encoders.

C. Imitation Learning with Invariant Representation

After learning the invariant representation, we use it for imitation in the GAIL fashion. The only difference from GAIL is that we use the encoded invariant states and actions in the calculation. So the loss function for the discriminator is defined as

$$\mathcal{L}_d = -\mathbb{E}_{\pi_E}[\log(D(z_t^s, z_t^a))] - \mathbb{E}_{\pi}[\log(1 - D(z_t^s, z_t^a))]. \quad (10)$$

The reward for the agent R_c is defined as $r_t = -\log(1 - D(z_t^s, z_t^a))$. The overall imitation learning process is summarized in Algorithm 1.

Algorithm 1: IR-GAIL

Input: Policy π and discriminator D . Agent configuration c

```
// Step 1: Representation Learning
1 Build training dataset by random rollout and expert demo.;
2 for  $i = 1, 2, \dots$  do
3   | Optimize  $\Phi_s, \Phi_a, \Gamma_s, \Gamma_a, \mathcal{F}$  by minimizing  $\mathcal{L}$ ;
4   | Optimize  $T_s, T_a$  by minimizing  $\mathcal{L}_T$ ;
5 end
// Step 2: Imitation Learning
6 for  $i = 1, 2, \dots$  do
7   | Sample trajectories using  $\pi_{\theta_i}$ ;
8   | Update  $D$  by minimizing  $\mathcal{L}_d$ ;
9   | Update  $\pi$  based on  $r_t$  using PPO [19];
10 end
```

V. EXPERIMENTS

A. Settings

We validate our algorithm on several MuJoCo control benchmarks. We use InvertedPendulum, InvertedDoublePen-

dulum, Hopper, Walker, Swimmer, and Halfcheetah. We set up the configuration space \mathcal{C} for these robots as follows. The values shown below are relative to the default configuration value in MuJoCo.

1) *InvertedPendulum*: The configuration space is $[0.75, 5.0] \times [0.5, 2.0]$. The first dimension corresponds to the length of the link. The second dimension corresponds to the maximum gear.

2) *InvertedDoublePendulum*: The configuration space is $[1.0, 3.0] \times [1.0, 3.0] \times [0.5, 1.5]$. The first and the second dimensions are the length of the bottom link and the above link respectively. The third dimension is the weight of the above link.

3) *Hopper*: The configuration space is $[0.5, 3.0] \times [0.5, 3.0]$. The first and the second dimension correspond to the length of the body and the length of the thigh respectively.

4) *Walker*: The configuration space is $[0.5, 2.0] \times [0.5, 2.0]$. The first and the second dimension correspond to the length of the thigh and the shank respectively.

5) *Swimmer*: The configuration space is $[0.5, 2.0] \times [0.5, 2.0] \times [0.5, 2.0]$. The first, second, and third dimension correspond to the length of the head, body, and the tail respectively.

6) *Cheetah*: The configuration space is $[0.5, 2.5] \times [0.5, 2.0] \times [0.5, 2.0]$. The first and the second dimension correspond to the length of the body and back leg respectively. The third dimension corresponds to the maximum gear of the back leg.

We use the default sensor data as state observation in the experiments. For the InvertedPendulum and Cheetah, we use two different types of state observation. The first type is based on the body keypoints (K). The second type is based on the joint angles (A), which is the same as the default MuJoCo setup.

B. Evaluation

We evaluate the performance of algorithms by interpolation and extrapolation experiments. In the interpolation experiment, we sample robots whose physical configurations are quite close to the experts'. In the extrapolation experiment, however, the configurations of sampled robots differ significantly from the experts'. For each environment, the expert demonstrations and random rollout used for training are collected on robots sampled from a ball region \mathcal{B} in the configuration space, so that we can sample robots in (out of) \mathcal{B} for interpolation (extrapolation) evaluation. We collect the

TABLE I
THE INTERPOLATION (INT.) AND EXTRAPOLATION (EXT.) PERFORMANCE OF EVALUATED ALGORITHMS.

Mode	Algorithm	IPendulum (K)	IDPendulum	Swimmer	Hopper	Walker	Cheetah (K)	IPendulum (A)	Cheetah (A)
Int.	GAIL	1.00±0.00	0.82±0.26	0.72±0.10	0.90±0.02	0.58±0.04	0.49±0.25	1.00±0.00	0.76±0.09
	TPIL	1.00±0.00	0.85±0.13	0.86±0.05	0.93±0.02	0.64±0.04	0.58±0.21	1.00±0.00	0.85±0.04
	IR-GAIL	1.00±0.00	0.96±0.04	0.98±0.02	0.99±0.01	0.83±0.08	0.91±0.04	1.00±0.00	0.86±0.06
Ext.	GAIL	0.65±0.47	0.04±0.01	0.65±0.11	0.78±0.05	0.45±0.12	0.22±0.14	1.00±0.00	0.71±0.10
	TPIL	0.82±0.38	0.05±0.01	0.77±0.06	0.80±0.04	0.57±0.06	0.46±0.15	1.00±0.00	0.81±0.08
	IR-GAIL	1.00±0.00	0.72±0.18	0.91±0.03	0.92±0.02	0.71±0.06	0.80±0.21	1.00±0.00	0.79±0.11

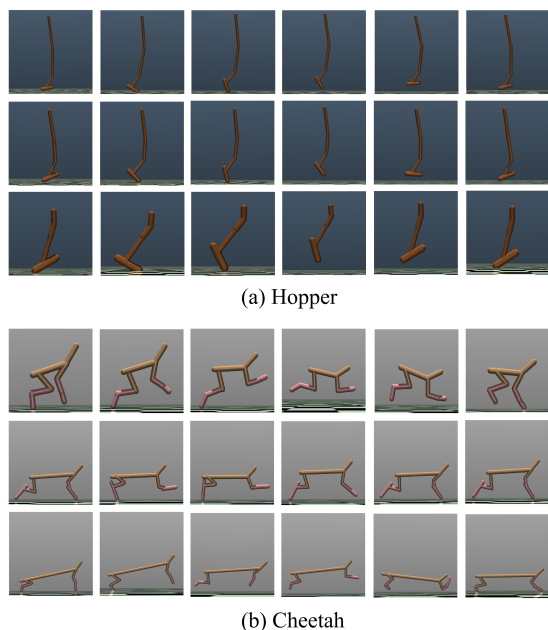


Fig. 5. Qualitative coupling results. Each row corresponds to states of one robot instance. Each column shows a group of coupled states. The coupled states in the same group are functionally similar.

random rollout on 16 randomly sampled robots. We sample 32 trajectories from 4 different experts for imitation. For policy optimization, we use PPO with Generalized Advantage Estimation [18].

C. Results

We compare the proposed IR-GAIL with GAIL, and another domain confusion based CDIL algorithm called TPIL [22]. The input of TPIL is fully observational, so we have to adapt it to our setting. The result is shown in Table 1. We calculate the mean and standard deviation of each algorithm’s return, which are normalized with respect to the average return of the expert policy and the random policy. We can see that our proposed IR-GAIL in general performs better than GAIL and TPIL, though these methods can still learn from the expert on some of these domains, which is due to the similarity of observations and dynamics. IR-GAIL can also generalize well to robots out of the training distribution, though the extrapolation performance of IR-GAIL is lower than interpolation performance. We also notice that the previous algorithms can still perform quite well on InvertedPendulum and Cheetah, if the observation is based on the joint angle. Such joint angle based observation is naturally invariant representations in some cases. For example, if the robots are only different from each other in size, then the angle information is still useful among them. However, measuring the angle accurately in the application can be difficult, and keypoint-based observation is used more frequently.

D. Ablation Studies

We also test the performance of the model after removing the dynamics regularization. We find that removing such regularization will not lead to a performance drop on InvertedPendulum and Hopper. However, we observe the performance on Walker, Swimmer, Cheetah, and InvertedDoublePendulum drops by 12%, 19%, 26%, and 43% respectively on average after its removal. An interesting fact is that the performance after such removal is still higher than the baselines in general. The reason is that variation bottleneck and the disentangle loss naturally regularize the representation and make it robust, which is also reported by [3].

E. Qualitative Results

In this part, we provide some qualitative coupling results to understand the learned invariant representation space in Figure 5. We first rollout trajectories of different robot instances by trained policies, and map the encountered states into the invariant representation space using the trained invariant representation module. To obtain a group of coupled states of different robot instances, we sample a point in the representation space and find its nearest representations of different robot instances. Then the corresponding robot states of these representations form a group of coupled states. The results are collected from Hopper and Cheetah. We can find that the coupled states in each group are highly functionally similar. This result suggests that our method can encode desired invariant representation for similar robots.

VI. CONCLUSION

In this paper, we studied the CDIL problem on a class of similar robots, and propose a new CDIL algorithm based on invariant representation. Experimental results showed that our method can achieve superior performance. We used visualization to demonstrate that our method can learn invariant representation for CDIL. One limitation of this work is that it is not observational, and we will study how to extend this approach to observational IL in the future. Another future direction is to study how to infer a robot’s configuration automatically. We will also try to apply the proposed method to real robots.

REFERENCES

- [1] M. Bain and C. Sammut. A framework for behavioural cloning. In *Machine Intelligence 15*, pages 103–129, 1995.
- [2] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm. Mutual information neural estimation. In *International Conference on Machine Learning (ICML)*, pages 531–540, 2018.
- [3] E. Cetin and O. Celiktutan. Domain-robust visual imitation learning with mutual information constraints. In *International Conference on Learning Representations (ICLR)*, 2021.
- [4] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- [5] A. Gupta, C. Devin, Y. Liu, P. Abbeel, and S. Levine. Learning invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint arXiv:1703.02949*, 2017.

- [6] D. Hejna, L. Pinto, and P. Abbeel. Hierarchically decoupled imitation for morphological transfer. In *International Conference on Machine Learning (ICML)*, pages 4159–4171. PMLR, 2020.
- [7] J. Ho and S. Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems (NIPS)*, volume 29, pages 4565–4573, 2016.
- [8] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- [9] K. Kim, Y. Gu, J. Song, S. Zhao, and S. Ermon. Domain adaptive imitation learning. In *International Conference on Machine Learning (ICML)*, pages 5286–5295. PMLR, 2020.
- [10] Y. Liu, A. Gupta, P. Abbeel, and S. Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.
- [11] Y. Lu and J. Tompson. Adail: Adaptive adversarial imitation learning. *arXiv preprint arXiv:2008.12647*, 2020.
- [12] A. N. Meltzoff and P. J. Marshall. Importance of body representations in social-cognitive development: new insights from infant brain science. *Progress in brain research*, 254:25–48, 2020.
- [13] A. Y. Ng, S. J. Russell, et al. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning (ICML)*, volume 1, page 2, 2000.
- [14] X. B. Peng, E. Coumans, T. Zhang, T. E. Lee, J. Tan, and S. Levine. Learning agile robotic locomotion skills by imitating animals. In *Robotics: Science and Systems (RSS)*, 2020.
- [15] X. B. Peng, A. Kanazawa, S. Toyer, P. Abbeel, and S. Levine. Variational discriminator bottleneck: Improving imitation learning, inverse rl, and gans by constraining information flow. In *International Conference on Learning Representations (ICLR)*, 2019.
- [16] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- [17] K. Schmeckpeper, O. Rybkin, K. Daniilidis, S. Levine, and C. Finn. Reinforcement learning with videos: Combining offline observations with interaction. In *Conference on Robot Learning (CoRL)*, 2020.
- [18] J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. In *International Conference on Learning Representations (ICLR)*, 2016.
- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [20] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine. AVID: learning multi-stage tasks via pixel-level translation of human videos. In *Robotics: Science and Systems (RSS)*, 2020.
- [21] Y. Song, L. Yu, Z. Cao, Z. Zhou, J. Shen, S. Shao, W. Zhang, and Y. Yu. Improving unsupervised domain adaptation with variational information bottleneck. *arXiv preprint arXiv:1911.09310*, 2019.
- [22] B. C. Stadie, P. Abbeel, and I. Sutskever. Third person imitation learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- [23] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [24] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30. IEEE, 2017.
- [25] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5026–5033, 2012.
- [26] F. Torabi, G. Warnell, and P. Stone. Generative adversarial imitation from observation. *arXiv preprint arXiv:1807.06158*, 2018.
- [27] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- [28] T. R. Zentall. Imitation in animals: evidence, function, and mechanisms. *Cybernetics & Systems*, 32(1-2):53–96, 2001.
- [29] Q. Zhang, T. Xiao, A. A. Efros, L. Pinto, and X. Wang. Learning cross-domain correspondence for control with dynamics cycle-consistency. In *International Conference on Learning Representations (ICLR)*, 2021.
- [30] Z. Zhu, K. Lin, and J. Zhou. Transfer learning in deep reinforcement learning: A survey. *arXiv preprint arXiv:2009.07888*, 2020.