# Learning Hierarchical Policies from Unsegmented Demonstrations using Causal Information

Arjun Sharma\*, Mohit Sharma\*, Nick Rhinehart, Kris M. Kitani Robotics Institute, Carnegie Mellon University

Abstract—The use of imitation learning to learn a single policy for a complex task that has multiple modes or hierarchical structure can be challenging. In fact, previous work has shown that learning separate policies for each mode or sub-task can greatly improve the performance of imitation learning. In this work, we model the interaction between sub-tasks and their resulting state-action trajectory sequences as a directed graphical model. We propose a new algorithm based on the generative adversarial imitation learning framework which learns subtask policies from unsegmented demonstrations. Our approach maximizes the causal information flow in the graphical model between sub-task latent variables and their generated trajectories. We also show how our approach connects with existing 'Options' framework commonly used to learn hierarchical policies.

## I. INTRODUCTION

Complex human activities can often be broken down into various simpler sub-activities or sub-tasks that can serve as the basic building blocks for completing a variety of complicated tasks. For instance, when driving a car, a driver may perform several simpler sub-tasks such as driving straight in a lane, changing lanes, executing a turn and braking, in different orders and for varying times depending on the source, destination, traffic conditions etc. Using imitation learning to learn one monolithic policy for each activity can be challenging as it ignores the shared sub-structure among the various activities. In this work, we develop an imitation learning framework that can learn a policy for each of these sub-tasks given unsegmented activity demonstrations and also learn a macro-policy which dictates switching from one sub-task policy to another. Learning sub-task specific policies has the benefit of shared learning. Each such sub-task policy also needs to specialize over a restricted state space, thus making the learning problem easier.

Previous works in imitation learning [6, 3] focus on learning each sub-task specific policy using *segmented* expert demonstrations by modeling the variability in each sub-task policy using a latent variable. This latent variable is inferred by enforcing high mutual information between the latent variable and expert demonstrations. This information theoretic perspective is equivalent to the graphical model shown in Figure 1 (Left), where the node c represents the latent variable. However, since learning sub-task policies requires isolated demonstrations for each sub-task this setup is difficult to scale to many real world scenarios where providing such segmented trajectories is cumbersome. Further, this setup does not learn a macro-policy to combine the learned sub-task policies in meaningful ways to achieve different tasks.

In our work we aim to learn each sub-task policy directly



Fig. 1: Left: Graphical model used in Info-GAIL [6]. Right: Causal model in this work. The latent code causes the policy to produce a trajectory. The current trajectory, and latent code produce the next latent code

from *unsegmented* activity demonstrations. For example, given a task consisting of three sub-tasks — A, B and C, we wish to learn a policy to complete sub-task A, learn when to transition from A to B, finish sub-task B and so on. To achieve this we use a causal graphical model, which can be represented as a Bayesian Network as shown in Figure 1 (Right). The nodes  $c_t$  denote latent variables which indicate the currently active sub-task and the nodes  $\tau_t$  denote the state-action pair at time t. We consider as given, a set of expert demonstrations, each of which is represented by  $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_T\}$  and has a corresponding sequence of latent factors  $\boldsymbol{c} = \{c_1, \dots, c_{T-1}\}$ . The sub-activity at time t dictates what state-action pair was generated at time t. The previous sub-task and the current stateaction pair together cause the selection of the next sub-task.

As we discuss in Section III, extending the use of mutual information to learn sub-task policies from unsegmented demonstrations is problematic, as it requires learning the macro-policy as a conditional probability distribution which depends on the unobserved future. This unobserved future is unknown during earlier points of interaction (Figure 1). To alleviate this, in our work we aim to force the policy to generate trajectories that maximize the directed information or causal information [7] flow from trajectories to latent factors of variation within the trajectories instead of mutual information. Using causal information requires us to learn a causally conditioned probability distribution [5] which depends only on the observed past while allowing the unobserved future to be sequentially revealed. Further, since there exists feedback in our causal graphical model *i.e.*, information flows from the latent variables to trajectories and vice versa, causal information also provides a better upper bound on this information flow between the latent variables and expert trajectories than does the conventional mutual information [7, 5].

We also draw connections with existing work on learning sub-task policies using imitation learning with the *options framework* [9, 2]. We show that our work, while derived using the information theoretic perspective of maximizing causal information, bears a close resemblance to applying the options framework in a generative adversarial imitation setting. Thus, our approach combines the benefits of learning hierarchical policies using the options framework with the robustness of generative adversarial imitation learning, helping overcome problems such as compounding errors that plague behaviour cloning.

In summary, the main contributions of our work include:

- We extend existing generative adversarial imitation learning frameworks to allow for learning of sub-task specific policies by maximizing causal information in a causal graph of sub-activity latent variables and observed trajectory variables.
- We draw connections between previous works on imitation learning with sub-task policies using *options* and show that our proposed approach resembles learning of *options* in a generative adversarial setting.

## II. RELATED WORK

# A. Imitation Learning

Imitation Learning [8] aims at learning policies that can mimic expert behaviours from demonstrations. Modeling the problem as a Markov Decision Process (MDP), the goal in imitation learning is to learn a policy  $\pi(a|s)$ , which defines the conditional distribution over actions  $a \in \mathcal{A}$  given the state  $s \in \mathcal{S}$ , from state-action trajectories  $\tau = (s_0, a_0, \cdots, s_T)$ of expert behaviour. Recently, [4] introduced an imitation learning framework called Generative Adversarial Imitation Learning (GAIL) that is able to learn policies for complex high-dimensional physics-based control tasks. The reduce the imitation learning problem into an adversarial learning framework, for which they utilize Generative Adversarial Networks (GAN). The generator network of the GAN represents the agent's policy  $\pi$  while the discriminator network serves as a local reward function and learns to differentiate between state-action pairs from the expert policy  $\pi_{\mathbb{E}}$  and from the agent's policy  $\pi$ . Mathematically, it is equivalent to solving the following optimization problem,

$$\min_{\pi} \max_{D} \mathbb{E}_{\pi}[\log D(s, a)] + \mathbb{E}_{\pi_{E}}[1 - \log D(s, a)] - \lambda H(\pi)$$

InfoGAIL [6] and [3] solve the problem of learning from policies generated by a mixture of experts. They introduce a latent variable c into the policy function  $\pi(a|s, c)$  to separate different type of behaviours present in the demonstration. To incentivize the network to use the latent variable, they utilize an information-theoretic regularization enforcing that there should be high mutual information between c and the state-action pairs in the generated trajectory, a concept that was first introduced in InfoGAN [1]. They introduce a variational lower bound  $L_1(\pi, Q)$  of the mutual information  $I(c; \tau)$  to the loss function in GAIL.

$$L_1(\pi, Q) = \mathbb{E}_{c \sim p(c), a \sim \pi(\cdot|s, c)} \log Q(c|\tau) + H(c) \le I(c; \tau)$$

The modified objective can then be given as,

$$\min_{\pi,q} \max_{D} \mathbb{E}_{\pi}[\log D(s,a)] + \mathbb{E}_{\pi_{E}}[1 - \log D(s,a)] -\lambda_{1}L_{1}(\pi,q) - \lambda_{2}H(\pi)$$

InfoGAIL models variations between different trajectories as the latent codes correspond to trajectories coming from different demonstrators. In contrast, we aim to model *intra-trajectory variations* and latent codes in our work correspond to sub-tasks (variations) within a demonstration. In section III, we discuss why using a mutual information based loss is in-feasible in our problem setting and describe our proposed approach.

## B. Options

Consider an MDP with states  $s \in S$  and actions  $a \in A$ . Under the options framework [9], an option, indexed by  $o \in O$ consists of a sub-policy  $\pi(a|s, o)$ , a termination policy  $\pi(b|s, \bar{o})$ and an option activation policy  $\pi(o|s)$ . After an option is initiated, actions are generated by the sub-policy until the option is terminated and a new option is selected. In [2] the authors formulate the options framework as a probabilistic graphical model where options are treated as latent variables which are then learned from expert data. The option policies  $(\pi(a|s, o))$  are analogous to sub-task policies in our work, we connect our work to this existing method and show that our method can be seen as a generative adversarial variant of the approach in [2].

# III. PROPOSED APPROACH

As mentioned in the previous section, while prior approaches can learn to disambiguate the multiple modalities in the demonstration of a sub-task and learn to imitate them, they cannot learn to imitate demonstrations of unsegmented long tasks that are formed by a combination of many small sub-tasks. To learn such sub-task policies from unsegmented gestures we use the graphical model in Figure 1 (Right), *i.e.*, consider a set of expert demonstrations, each of which is represented by  $\tau = {\tau_1, \dots, \tau_T}$  where  $\tau_t$  is the state-action pair observed at time t. Each such demonstration has a corresponding sequence of latent variables  $c = {c_1, \dots, c_{T-1}}$  which denote the subactivity in the demonstration at any given time step.

As noted before, previous approaches [6, 3] model the expert sub-task demonstrations using only a single latent variable. To enforce the model to use this latent variable, previous approaches propose to maximize the mutual information between the demonstrated sequence of state-action pairs and the latent embedding of the nature of the sub-activity. This is achieved by adding a lower bound to the mutual information between the latent variables and expert demonstrations. This variational lower bound of the mutual information is then combined with the the adversarial loss for imitation learning proposed in [4]. Extending this to our setting, where we have a *sequence* of latent variables c, yields the following lower bound on the mutual information,

$$L(\pi, q) = \sum_{t} \mathbb{E}_{c^{1:t} \sim p(c^{1:t}), a^{t-1} \sim \pi(\cdot | s^{t-1}, c^{1:t-1})} \begin{bmatrix} \\ \log q(c^{t} | c^{1:t-1}, \tau) \end{bmatrix} + H(c) \le I(\tau; c)$$
(1)

Observe that the dependence of q on the entire trajectory  $\tau$  precludes the use of such a distribution at test time, where only the trajectory up to the current time is known. To overcome this limitation, in this work we propose to force the policy to generate trajectories that maximize the *directed* or *causal* information flow from trajectories to the sequence of latent sub-activity variables instead. As we show below, by using causal information instead of mutual information, we can replace the dependence on  $\tau$  with a dependence on the trajectory generated up to current time t.

The causal information flow from a sequence X to Y is given by,

$$I(\mathbf{X} \to \mathbf{Y}) = H(\mathbf{Y}) - H(\mathbf{Y} \| \mathbf{X})$$

where H(Y||X) is the causally-conditioned entropy. Replacing X and Y with the sequences  $\tau$  and c gives,

$$I(\boldsymbol{\tau} \to \boldsymbol{c}) = H(\boldsymbol{c}) - H(\boldsymbol{c} \| \boldsymbol{\tau})$$
  
=  $H(\boldsymbol{c}) - \sum_{t} H(c^{t} | c^{1:t-1}, \tau^{1:t})$   
=  $H(\boldsymbol{c}) + \sum_{t} \sum_{c^{1:t-1}, \tau^{1:t}} \left[ p(c^{1:t-1}, \tau^{1:t}) \right]$   
 $\sum_{c^{t}} p(c^{t} | c^{1:t-1}, \tau^{1:t}) \log p(c^{t} | c^{1:t-1}, \tau^{1:t}) \right]$   
(2)

A variational lower bound,  $L_1(\pi, q)$  of the causal information,  $I(\tau \rightarrow c)$  which uses an approximate posterior  $q(c^t|c^{1:t-1}, \tau^{1:t})$  instead of the true posterior  $p(c^t|c^{1:t-1}, \tau^{1:t})$  can then be derived to get,

$$L_{1}(\pi, q) = \sum_{t} \mathbb{E}_{c^{1:t} \sim p(c^{1:t}), a^{t-1} \sim \pi(\cdot | s^{t-1}, c^{1:t-1})} \begin{bmatrix} \\ \log q(c^{t} | c^{1:t-1}, \tau^{1:t}) \end{bmatrix} + H(\mathbf{c}) \leq I(\mathbf{\tau} \to \mathbf{c})$$
(3)

Thus, by maximizing causal information instead of mutual information, we can learn a posterior distribution over the next latent factor c given the latent factors discovered up to now and the trajectory followed up to now, thereby removing the dependence on the future trajectory. In practice, we do not consider the H(c) term. This gives us the objective,

$$\min_{\pi,q} \max_{D} \quad \mathbb{E}_{\pi}[\log D(s,a)] + \mathbb{E}_{\pi_{E}}[1 - \log D(s,a)] - \lambda_{1}L_{1}(\pi,q) - \lambda_{2}H(\pi)$$
(4)

We call this approach Causal-Info GAIL. Notice, that to compute the loss in equation 3, we need to sample from the



Fig. 2: VAE pre-training step. The VAE encoder uses the current state  $(s_t)$ , and previous latent variable  $(c_{t-1})$  to produce the current latent variable  $(c_t)$ . The decoder reconstructs the action  $(a_t)$  using  $s_t$  and  $c_t$ .

prior distribution  $p(c^{1:t})$ . In order to estimate this distribution, we train a variational auto-encoder (VAE) on the expert trajectories. Figure 2 shows the design of the VAE pictorially. We use the following objective, which maximizes the lower bound of the probability of the trajectories  $p(\tau)$ , to train our VAE,

$$L_{\text{VAE}}(\pi, q; \boldsymbol{\tau}_{i}) = -\sum_{t} \mathbb{E}_{c^{t} \sim q} \Big[ \log \pi(a^{t} | s^{t}, c^{1:t}) \Big] + \sum_{t} D_{\text{KL}}(q(c^{t} | c^{1:t-1}, \tau^{1:t}) \| p(c^{t} | c^{1:t-1}))$$
(5)

We can then use q to obtain samples of latent variable sequence c by using the expert demonstrations.

# A. Connection with options framework

In [2] the authors provide a probabilistic perspective of the options framework. Although, [2] consider separate termination and option latent variables ( $b^t$  and  $o^t$ ), for the purpose of comparison, we collapse them into a single latent variable  $c^t$ , similar to our framework with a distribution  $p(c^t|s^t, c^{t-1})$ . The lower-bound derived in [2] which is maximized using EM can then be written as (suppressing dependence on parameters),

$$p(\tau) \ge \sum_{t} \sum_{c^{t-1:t}} p(c^{t-1:t}|\tau) \log p(c^{t}|s^{t}, c^{t-1})) + \sum_{t} \sum_{c^{t}} p(c^{t}|\tau) \log \pi(a^{t}|s^{t}, c^{t})$$
(6)

Note that the first term in Equation (6) *i.e.*, the expectation over the distribution  $\log p(c^t | s^t, c^{t-1})$  is the same as Equation (3) of our proposed approach with a one-step Markov assumption and a conditional expectation with given expert trajectories instead of an expectation with generated trajectories. The second term in Equation (6) *i.e.*, the expectation over  $\log \pi(a^t | s^t, c^t)$  is replaced by the GAIL loss in Equation (4). Our proposed Causal-Info GAIL can be therefore be considered as the generative adversarial variant of imitation learning using the options framework. The VAE behaviour cloning pre-training step in Equation (5) is exactly equivalent to Equation (6), where we use variational inference instead of EM. Thus, our approach combines the benefits of both behavior cloning and generative



Fig. 3: Results on the Four Rooms environment. Each figure shows results for a different latent variable. The arrows in each cell indicate the direction (action) with highest probability in that state and using the given latent variable.

adversarial imitation learning. Using GAIL enables learning of robust policies that do not suffer from the problem of compounding errors. At the same time, conditioning GAIL on latent codes learned from the behavior cloning step prevents the issue of mode collapse in GANs.

## **IV. EXPERIMENTS**

## A. Environment

We test our approach on a grid world environment. The environment is a  $15 \times 11$  grid and consists of four rooms connected via corridors as shown in Figures 3 and 4. The agent spawns at a random location in the grid at the beginning of each episode. One of the four rooms is then selected at random and an apple is placed at the centre of the room. The goal of the agent is to find the shortest path to the apple using four actions — up, down, left and right. Expert trajectories were created using Dijkstra's shortest path algorithm.

### B. Results

Figure 3 shows sub-task policies learned by our approach. Each of the four plots corresponds to a different value of the latent variable. The arrow at every state in the grid in each plot shows the agent action (direction) with the highest probability in that state for that latent variable. In the discussion that follows, we label the rows from 1 to 4 starting from the room at the top left and moving in the clockwise direction.

Notice how the different latent variables correspond to different sub-tasks. For e.g., the latent code in Figure 3(b) learns the sub-task of moving from room 1 to room 3 and from room 2 to room 4 and the code in Figure 3(c) learns the



Fig. 4: Expert and generated trajectories in the Four Rooms environment. Star (\*) represents the start state. The expert trajectory is shown in red. The color of the generated trajectory represents the latent code used by the policy at each time step.

opposite sub-tasks. The code in Figure 3(d) learns the sub-task of moving from rooms 2 and 4 to the horizontal corridor.

Figure 4 shows two examples of how the macro-policy switches between various latent codes to achieve the desired goals of reaching the apples in rooms 1 and 2 respectively.

#### REFERENCES

- Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2172–2180, 2016.
- [2] Christian Daniel, Herke Van Hoof, Jan Peters, and Gerhard Neumann. Probabilistic inference for determining options in reinforcement learning. *Machine Learning*, 104(2-3): 337–357, 2016.
- [3] Karol Hausman, Yevgen Chebotar, Stefan Schaal, Gaurav Sukhatme, and Joseph J Lim. Multi-modal imitation learning from unstructured demonstrations using generative adversarial nets. In Advances in Neural Information Processing Systems, pages 1235–1245, 2017.
- [4] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In Advances in Neural Information Processing Systems, pages 4565–4573, 2016.
- [5] Gerhard Kramer. Directed information for channels with feedback. PhD thesis, Eidgenossiche Technische Hochschule Zurich, 1998.
- [6] Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual demonstrations. In Advances in Neural Information Processing Systems, pages 3815–3825, 2017.
- [7] James Massey. Causality, feedback and directed information. In Proc. Int. Symp. Inf. Theory Applic.(ISITA-90), pages 303–305. Citeseer, 1990.
- [8] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In Advances in neural information processing systems, pages 305–313, 1989.
- [9] Richard S Sutton, Doina Precup, and Satinder P Singh. Intra-option learning about temporally abstract actions. In *ICML*, volume 98, pages 556–564, 1998.