

# Case-Based Inverse Reinforcement Learning Using Temporal Coherence

Jonas Nüßlein \*\*, Steffen Illium, Robert Müller, Thomas Gabor, and  
Claudia Linnhoff-Popien

LMU Munich, Institute of Computer Science, Germany,  
{jonas.nuesslein, steffen.illum, robert.mueller, thomas.gabor,  
linnhoff}@ifi.lmu.de

**Abstract.** Providing expert trajectories in the context of Imitation Learning is often expensive and time-consuming. The goal must therefore be to create algorithms which require as little expert data as possible. In this paper we present an algorithm that imitates the higher-level strategy of the expert rather than just imitating the expert on action level, which we hypothesize requires less expert data and makes training more stable. As a prior, we assume that the higher-level strategy is to reach an unknown target state area, which we hypothesize is a valid prior for many domains in Reinforcement Learning. The target state area is unknown, but since the expert has demonstrated how to reach it, the agent tries to reach states similar to the expert. Building on the idea of Temporal Coherence, our algorithm trains a neural network to predict whether two states are similar, in the sense that they may occur close in time. During inference, the agent compares its current state with expert states from a Case Base for similarity. The results show that our approach can still learn a near-optimal policy in settings with very little expert data, where algorithms that try to imitate the expert at the action level can no longer do so.

**Keywords:** Case-Based Reasoning, Inverse Reinforcement Learning, Incomplete Trajectories, Learning from Observations, Temporal Coherence

## 1 Introduction

In Reinforcement Learning (RL), the goal of the agent, given a Markov Decision Process, is to maximize the expected cumulative reward. The higher the expected reward, the better the agent's policy. In Imitation Learning, on the other hand, the agent does not have access to a reward signal from the environment. Instead, it either has access to an expert who can be asked online for the best action for a given state or a set of trajectories generated by the expert is available. Imitation Learning has been proven to be particularly successful in domains where the demonstration by an expert is easier than the construction of a suitable reward function [AD21]. There are two main approaches to Imitation Learning:

---

\*\* Corresponding author

Behavioral Cloning (BC) [Pom91] and Inverse Reinforcement Learning (IRL) [FLL17,AD21,AN04]. In BC, the agent learns via supervised learning to produce the same actions that the expert would have produced. The advantage of this approach is that no further rollouts in the environment are necessary. However, the approach suffers greatly from compounding error, i.e., the slow drift of states visited by the expert [RGB11]. In the second approach, IRL, a reward function is learned under which the expert is uniquely optimal. Then, a policy can be learned using classical Reinforcement Learning and this reconstructed reward function. However, the drawback of this approach is that it usually requires a lot of rollouts in the environment, as it often includes RL as a subroutine.

GAIL [HE16] is another approach to Imitation Learning. It builds on the ideas of Generative Adversarial Networks. In this approach, a policy and a discriminator are learned. The goal of the discriminator is to be able to distinguish state-action pairs of the expert from state-action pairs of the agent, while the goal of the policy is to fool the discriminator. GAIL requires expert actions, but there is an extension, named GAIFO, which does not [TWS18]. While GAIL discriminates between state-action pairs produced by the agent or the expert respectively, GAIFO does so with state transitions. In this paper we consider, as GAIFO does, the setting where no expert actions are available to the agent. This setting is also called Learning from Observation (LfO) or Imitation from Observation (IfO) [YMH<sup>+</sup>19,TWS18].

Providing expert trajectories is often very expensive and time-consuming, especially if the expert is a human. The goal must therefore be to create algorithms which require as little expert data as possible.

The aim of this paper is to present an algorithm that imitates the higher-level strategy of the expert rather than just imitating the expert on action level.

Our motivation for this is that we hypothesize that it takes less expert data to learn the higher-level strategy than to imitate the expert on action level. We also hypothesize that it makes the training more stable, with less “forgetting” of what has already been learned. As a prior for the higher-level strategy, we assume that the higher-level strategy is to reach an unknown target state area, which we hypothesize is a valid prior for many domains in Reinforcement Learning.

We present an algorithm that learns these higher-level strategies from expert trajectories. To prove that the algorithm does not imitate the expert on action level, we consider a special setting of Imitation Learning, which is characterized by incomplete expert trajectories. Here, the agent does not see every state of the expert trajectory, but, for example, only every fifth. Thus, it cannot imitate the expert on action level.

The idea behind machine learning is to derive general rules from a large amount of data, which can then be applied to new, unknown scenarios. This induction-based learning principle differs from Case-Based problem solving. In Case-Based Reasoning, a set of problems solved in the past is stored in a database. If a new, unknown problem is to be solved, the problem most similar to the current

situation is retrieved from the database and used to solve the current problem. Applications of Case-Based Reasoning range from explaining neural network outputs [LLCR18, KK19] over financial risk detection [LPS21] to medical diagnosis [CB16]. In our algorithm we build on ideas from Case-Based Reasoning as well as on the idea of Temporal Coherence.

Temporal Coherence [GBT<sup>+</sup>15, MCW09, ZNY11] originates from Video Representation Learning, where the idea is that two images, which occur shortly after each other in a video, are very likely to show the same object or objects. The two images should therefore have a similar representation. On the other hand, distant images should have different representations. The combination of this convergence and divergence, also called contrastive learning, can be used as a self-supervised training signal to learn semantically meaningful representations [KHW<sup>+</sup>21].

Our contribution in this paper is twofold. First, we propose the setting with incomplete expert trajectories without expert actions as a way to prove that the agent really learns the expert’s strategy and does not imitate the expert on action level. The prior we are using for the higher-level strategy is to reach an unknown target state area. Second, we present an algorithm that can learn such higher-level strategies and we test it on four typical domains of RL. The results show that our approach can still learn a near-optimal policy in settings with very little expert data, where IRL algorithms that try to imitate the expert at the action level can no longer do so.

## 2 Background

In this section we want to provide a brief introduction to Markov Decision Processes (MDP) [ADBB17]. A MDP is a tuple  $(S, A, T, R, \gamma)$ .  $S$  is a set of states, combined with a distribution of starting states  $p(s_0)$ .  $A$  is a set of actions the agent can perform.  $T$  is the transition function of the environment which computes the next state  $s_{t+1}$  given a state  $s_t$  at time  $t$  and an action  $a_t$ :  $T(s_{t+1}|s_t, a_t)$ . The property of  $T$  that the computation of  $s_{t+1}$  depends only on the last state  $s_t$  and not on  $s_{\tau < t}$  is also called the Markov property, hence the name Markov Decision Process.  $r_t = R(s_t, a_t, s_{t+1})$  is a reward function and  $\gamma \in [0; 1]$  is a discount factor. If  $\gamma < 1$ , immediate rewards are preferred compared to later rewards. An agent acts in a MDP using its policy  $\pi$ . The policy is a function which outputs an action  $a$  given a state  $s$ :  $\pi(s) = a$ . MDPs are often episodic, which means that the agent acts for  $T$  steps, after which the environment is reset to a starting state. The goal of the agent in a MDP is to maximize the expected return by finding the policy

$$\pi^* = \operatorname{argmax}_{\pi} E[R|\pi] \tag{1}$$

where the return  $R$  is calculated via:

$$R = \sum_{t=0}^{T-1} \gamma^t r_{t+1} \tag{2}$$

### 3 Related Work

**Combination of Case-Based Reasoning (CBR) and Reinforcement Learning (RL):** In [BRLdM09] the authors use Case-Based Reasoning (CBR) in the setting of Heuristic Accelerated Reinforcement Learning, where a heuristic function assists in action selection to accelerate exploration in the state-action space. In [ALUHMA08], Case-Based Reasoning is used to efficiently switch between previously stored policies learned with classical RL. A similar approach is taken by [WW14]. Most Imitation Learning algorithms try to imitate the expert skill step-by-step. In [LHYL19], a hierarchical algorithm is presented where this goal is mitigated. Instead, a policy is learned that reaches sub-goals, which in turn are sampled by a meta-policy from the expert demonstrations.

**Temporal Coherence in Reinforcement Learning:** Some papers have already investigated the use of Temporal Coherence in the context of Reinforcement Learning. For example, in [FDH<sup>+</sup>19] it was proposed to learn an embedding for the inputs of the Markov Decision Process, such that the euclidean distance in the embedding space is proportional to the number of actions the current agent needs to get from one state to the other. The byproduct of this is a policy that can theoretically reach any previously seen state on demand. A similar idea is followed in the context of goal-conditioned RL: In [LSSL21] a proximity function  $f(s, g)$  is learned that outputs a scalar proportional to the distance of the state  $s$  to the goal  $g$ . The distance then serves as a dense reward signal for a classical RL agent. This is especially beneficial when the environment’s reward function is sparse.

In [SLC<sup>+</sup>18], a special setting is considered where multiple observations are available simultaneously, showing the same state from different perspectives. An embedding is then learned so that contemporaneous observations have the same embedding and temporally distant observations have different embeddings. Thus, a perspective-invariant representation is learned, which contains semantic information. That paper also considers the case where only one perspective is available. In this case, the embeddings of two nearby inputs should be as similar as possible and temporally distant inputs should be as dissimilar as possible. We build on this idea of Temporal Coherence, although we do not learn an embedding. [DTLS18] extends the idea of [SLC<sup>+</sup>18] to input sequences to contrast movements.

In [SRM<sup>+</sup>18], the concept of Reachability Networks is already introduced, i.e., a network that classifies whether two states can occur in short succession in a trajectory. This network is then used as a curiosity signal to guide exploration in sparse reward domains. We build on this concept, but use it differently. While in [SRM<sup>+</sup>18] the agent searches for dissimilar states, the goal of the agent in our approach is to reach similar states (compared to expert states).

**Curriculum via Expert Demonstrations:** As we will see in the next section, the reconstructed reward function in our approach can be interpreted as an implicit curriculum. A related approach, which creates an explicit curriculum

using expert demonstrations, is [DHW21]. In that paper the expert trajectory is divided into several sections and state resetting to expert states is used to increase the difficulty of reaching the goal state. The sector from which expert states are sampled for resetting is gradually pushed away from the goal as the curriculum progresses. A similar approach is [HAE<sup>+</sup>20], which again uses resetting to starting states of varying difficulty.

**Unsupervised Perceptual Rewards for Imitation Learning:** the closest work to ours is [SXL16]. In that paper the authors examine how to use pre-trained vision models to reconstruct a reward function from few human video demonstrations. They do so by first splitting the human demo videos into segments, then selecting features of a pre-trained model which best discriminate between the segments and then using a reward function, which is based on these selected features, to learn a policy via standard RL algorithms. The biggest difference to our algorithm is that [SXL16] reconstructs the reward function entirely before training in the RL domain. In contrast, we learn the reward function and the policy at the same time.

## 4 Case-Based Inverse Reinforcement Learning (CB-IRL)

In this work, we consider a special setting of Imitation Learning that is characterized by two main features. First, there are no expert actions available to the agent and, second, the expert trajectories are incomplete, i.e., from the original sequence of MDP states of the expert  $[s_0, s_1, s_2, \dots, s_T]$ , the agent only sees, for example, every fifth state:  $[s_0, s_5, s_{10}, \dots]$ . This makes it impossible for the agent to imitate the expert at the action level. Given such a setting, we now propose the algorithm Case-Based Inverse Reinforcement Learning (CB-IRL). The architecture of CB-IRL consists of the Case Base ( $C$ ) and two neural networks, the Equality Net ( $E$ ) and the Policy ( $\pi$ ), see Figure 1.  $C$  is filled with the expert trajectories.

The basic idea is that the agent should not act in every step exactly as the expert would do, but instead imitate the higher-level strategy of the expert. We chose the task of reaching a target state area as the prior for the higher-level strategy. For example, for the OpenAI domain ‘MountainCar’ the target state area are the states where the car is on top of the mountain. For the Atari game ‘Pong’ the target state area would be the states where the agent has 21 points. The agent does not know the target state area, but since the expert has demonstrated how to reach the target state area, CB-IRL trains the agent to reach similar states as the expert.

Two states are “similar” in the context of Reinforcement Learning if it takes only few steps (actions) to get from one state to the other. Other approaches [FDH<sup>+</sup>19,SLC<sup>+</sup>18,DTLS18] try to learn a state-embedding such that the euclidean distance of the representations is proportional to the number of steps needed to get from one state to the other. We take a different approach and

instead train a neural network that accepts two states  $s_1$  and  $s_2$  as input and outputs a scalar  $E(s_1, s_2) = d$ ;  $E : S \times S \rightarrow [0; 1]$  to classify whether  $s_2$  can be reached within *windowFrame* steps from  $s_1$ . Thus, this is a classification and not a regression. We believe a classification is easier and more stable to learn compared to a regression, since it suffers less from the “moving target” problem. For example if we would predict the number of steps which are required to go from one state to the other, the target of this supervised learning tasks is heavily based on the current performance of the agent. In contrast, for near/far classification, it does not matter if the states are, for example, 30 or 40 steps apart if *windowFrame* = 10. In both cases the state pair gets the target 0 for supervised learning, since it shall be classified as dissimilar.

---

**Algorithm 1: CB-IRL**

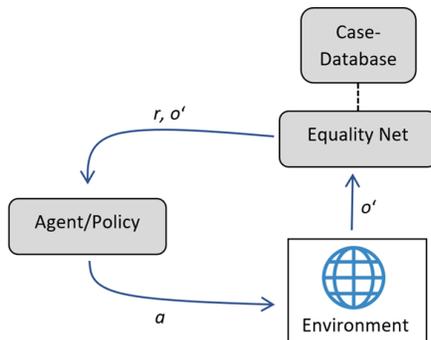

---

**Data:** Case-Base  $C$  (containing expert trajectories)  
**Result:** Policy  $\pi$ , Equality Net  $E$

**while** *training* **do**  
     $s \leftarrow$  sample start state  
     $r_{pre} \leftarrow$  Reward( $s$ )  
    *trajectory*  $\leftarrow [s]$   
    **while** *episode is not finished* **do**  
         $a \leftarrow \pi(s)$   
         $s' \leftarrow$  execute  $a$   
        *trajectory.append*( $s'$ )  
         $r_{post} \leftarrow$  Reward( $s'$ )  
         $r \leftarrow r_{post} - \alpha * r_{pre}$   
        use  $(s, a, r, s')$  for training  $\pi$   
         $s \leftarrow s'$   
         $r_{pre} \leftarrow r_{post}$   
    **end**  
    append *trajectory* to the Replay Buffer of  $E$   
    train  $E$  using the Replay Buffer,  $C$  and the hyperparameters  
    *windowFrame* and  $\nu$   
**end**

**Function** Reward( $s$ ):  
     $mostSimilar = \mu$   
     $similarity = \tau$   
    **foreach** *trajectory*  $\in C$  **do**  
        **foreach**  $o_e^{(i)} \in$  *trajectory* **do**  
            **if**  $E(s, o_e^{(i)}) > similarity$  **then**  
                 $mostSimilar = i$   
                 $similarity = E(s, o_e^{(i)})$   
            **end**  
        **end**  
    **end**  
    **return**  $mostSimilar$

---

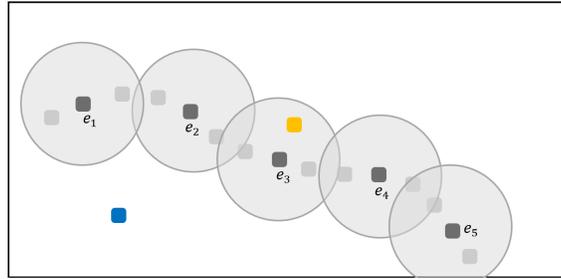


**Fig. 1.** This figure shows the usual cycle of Reinforcement Learning, with a small adjustment. The Equality Net ( $E$ ) is interposed between the environment and the agent. The agent performs an action  $a$ , which is executed in the environment.  $E$  receives the next observation  $o'$  from the environment, then calculates the reward  $r$  using the case database and forwards both to the agent.

A second advantage of Reachability Networks in contrast to embeddings is that they are suitable for asymmetric state-action spaces. For example, it may be easy to reach  $s_2$  from  $s_1$ , but difficult or impossible to reach  $s_1$  from  $s_2$ .

The policy  $\pi$  is learned via Inverse Reinforcement Learning using the case database  $C$  and the Equality Net  $E$ . If the agent is in state  $o$ , it executes the action  $a = \pi(o)$  with its current policy  $\pi$  and receives the next observation  $o'$  from the environment. Using  $E$ , all expert observations  $o_e^{(i)}$  from  $C$  are now checked to see if they are similar to  $o'$ , where the similarity must be above a threshold  $\tau$ . If there is a similar expert state  $o_e^{(j)}$  (if more than one, choose the most similar), the reward is given by the position number  $j$ . Thus, the further back the similar expert state is in the expert trajectory, the higher the reward the agent receives. If there is no similar expert state, the agent receives a penalty  $\mu$  (a negative reward). Figure 2 shows the idea schematically. The complete algorithm is summarized in Algorithm 1.

The algorithm contains several hyperparameters, whose task and influence we discuss in the following:  $\tau \in [0; 1]$  is the threshold that determines the minimum similarity of an expert state  $o_e^{(i)}$  to the current state  $o$  of the agent, so that the agent receives a positive reward. If no expert state has a similarity higher than  $\tau$ , the agent receives a penalty (a negative reward  $\mu$ ). The hyperparameter  $\alpha \in [0; 1]$  controls whether the actual reward for the agent is always the reward difference ( $\alpha = 1$ ) or whether the agent always receives the full reward ( $\alpha = 0$ ). For  $\alpha = 0$ , the agent tends to achieve large rewards as quickly as possible, but maybe not reliably, whereas for  $\alpha = 1$ , the agent tries to achieve a large reward as reliably as possible by the end of the episode.



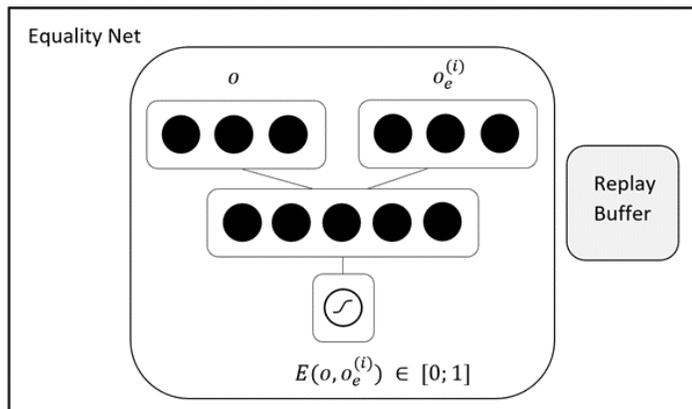
**Fig. 2.** This figure shows schematically how the algorithm works. Assume the rectangle is a two-dimensional state space. The light and dark gray boxes represent the trajectory of expert states but only the dark states are visible to the agent. Using its own rollouts the agent now learns the Equality Net, which classifies whether two states can occur close to each other in a trajectory. The outputs of the Equality Net for the expert states are represented in the image by the circles around them. During inference, the agent checks whether the current state is similar to an expert state or not. For example, if the agent is in the yellow state, it is similar to expert state  $e_3$  and therefore receives reward 3. If the agent is in the blue state, it is not similar to any expert state and receives a negative reward  $\mu$ .

The hyperparameters *windowFrame* and  $\nu$  are used to train  $E$ . They model on the one hand the threshold which indicates whether two states are considered similar or dissimilar and on the other hand the number of explicit divergence between states of the agent and states of the expert.

**Training of the Equality Net:** The task of the Equality Net  $E$  is to classify whether two inputs can occur in short succession in a trajectory and are thus “similar”. To train  $E$ , we use the Replay Buffer that contains the trajectories sampled by the agent.  $E$  is trained using supervised learning. The training set consists of similar and dissimilar state pairs. For the similar state pairs, two states are selected from the same trajectory of the Replay Buffer which are no further apart than *windowFrame* steps. For the dissimilar state pairs, two states are sampled from two different trajectories. For the similar state pairs, the network  $E$  is trained to output the value 1, for dissimilar state pairs it is trained to output 0. The structure of  $E$  is graphically visualized in Figure 3.

In addition, training can also be performed in an analogous manner on the expert trajectories. The hyperparameter  $\nu$  models the number of explicit divergence between agent and expert state. That is, there are  $\nu$  state pairs where one state is sampled from the Replay Buffer and the other state is sampled from  $C$ . The target for these state pairs during supervised learning is 0, since they shall be classified as dissimilar.

The output of the Equality Net can be understood as a (lossy) binary distance measure. The distance measure is binary because it only distinguishes between similar (1) and dissimilar (0).



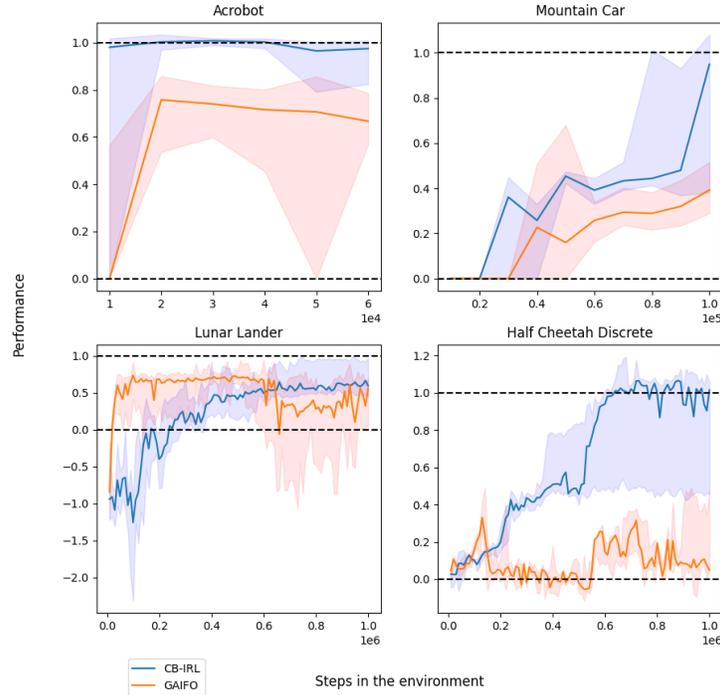
**Fig. 3.** The Equality Net  $E$  accepts as input two states and classifies whether they are similar in the sense that they can appear close to each other in a trajectory. During inference,  $E$  receives as input the current state  $o$  of the agent and compares it to all expert states  $o_e^{(i)}$ .  $E$  is trained using supervised learning on the trajectories produced by the agent, which are stored in the Replay Buffer.

## 5 Experiments

We tested our algorithm in four OpenAI Gym domains [BCP<sup>+</sup>16]: Acrobot, Mountain Car, Lunar Lander, and Half Cheetah. For Half Cheetah, we created a modified version called Half Cheetah Discrete. Details can be found in Appendix A. As justified in [CLB<sup>+</sup>17], only domains should be used for the evaluation of IRL algorithms in which the episodes are always of the same length. This is because early ending of episodes may contain implicit information about the reward. For example, in the ‘Mountain Car’ domain, the episode ends when the car has successfully driven up the hill. For this reason, we have adjusted all domains so that episodes are always of the same length, with the agent receiving the last observation until the end if the episode ended early.

We first trained an expert for each domain using the reward function of the environment. We then used these experts to create exactly one trajectory for each domain, which consisted only of the expert states and not the expert actions. We then used it to train CB-IRL and GAIFO. GAIFO had access to all expert states, while CB-IRL only had access to every tenth expert state. For example, for Lunar Lander, the expert trajectory was about 150 steps long, so the training set for GAIFO consisted of these 150 expert states, while the training set for CB-IRL consisted of only 15 expert states.

For the hyperparameter search, we tested five hyperparameter sets for each algorithm and domain and selected the best one. Using these hyperparameters, we then trained CB-IRL and GAIFO three times with three different seeds. During training we generated 20 episodes every 10,000 steps for each seed and



**Fig. 4.** Scaled performance of CB-IRL and GAIFO on four different domains trained using one expert trajectory, where GAIFO had access to all expert states and CB-IRL had access to only one in ten.

algorithm (for a total of 60 episodes per algorithm every 10,000 steps). For each episode, we calculated the total return using the environment’s reward function. The returns were then scaled using the performance of a random agent (representing value 0) and the expert (representing value 1). We then calculated the 0.25, 0.5 (median), and 0.75 quantiles for these 60 return values. For both algorithms, the solid lines represent the median and the shaded areas enclose the 0.25 and 0.75 quantiles.

As can be seen in Figure 4, CB-IRL mostly performed better than GAIFO in the experiments, even though it had access to only one tenth of GAIFO’s training set. Furthermore, CB-IRL showed a more stable learning behavior. The difference was particularly clear in the Half Cheetah Discrete domain. Here, the advantage of CB-IRL became apparent, where the agent did not learn to behave exactly like the expert in every state, but to reach similar states as the expert. CB-IRL has learnt the high-level strategy to “run as far as possible”.

A Python implementation of CB-IRL and the code used to create the experiments are available on GitHub [<https://github.com/JonasNuesslein/CB-IRL>]. For GAIFO we used the implementation of tf2rl [Ota20]. The chosen hyperparameters for the experiments can also be found on GitHub in the file *config.py*.

## 6 Discussion of the Approach

In this section we discuss the advantages and disadvantages of CB-IRL. Turning first to the disadvantages: The computation of the reward is more computationally intensive than in many other IRL algorithms, because in each step the current state must be compared against all expert states in the case base  $C$ . The run-time complexity is thus linear in the size of  $C$ . This can be serious for larger case bases, however, the target application areas of CB-IRL are precisely the settings where very little expert data is available. Moreover, the computational intensity can be reduced by calculating a reward only in every  $k$ -th step, rather than in every step.

The second drawback of our approach is the specialization of CB-IRL to state-reaching in contrast to state-keeping domains. By state-reaching domains, we mean domains in which certain variables of the state vector have to be changed. An example of this is the OpenAI Gym domain ‘Mountain Car’ [BCP<sup>+</sup>16], in which the goal is to maximize the x-position of the car. Another example is the Atari game ‘Pong’ [MKS<sup>+</sup>15], in which the goal of the agent is to reach 21 points. By state-keeping domains, we mean domains in which the goal is to leave certain variables of the state vector unchanged. An example of this would be ‘Cart-Pole’ [BCP<sup>+</sup>16], where the goal is to keep the angle of the pole at 90° if possible or ‘HalfCheetah’ [BCP<sup>+</sup>16], where the goal is to keep a high velocity. Due to the structure of CB-IRL, it is predominantly suitable for state-reaching domains, as the algorithm encourages the agent to reach states from the posterior of the expert trajectory.

The advantages of CB-IRL are that it does not require a reward function, expert actions, or complete expert trajectories. Since the agent can learn with incomplete expert trajectories, it has proven that it imitates the higher-level strategy of the expert and does not imitate the expert on action level.

This allows the agent to learn a near-optimal strategy with little data, which would be insufficient to imitate the expert on action level (as can be seen in the Half Cheetah Discrete domain). The learning behavior also shows a more stable pattern with less “forgetting” of what has already been learned.

A second possible advantage, which we leave as future work to verify, is that the Equality Net is not task-specific and can be reused for other tasks in the same domain, which can enable fast transfer learning.

A third possible advantage also left for future work is that the ability to learn from incomplete trajectories may be beneficial in real-world applications, where state observations may be noisy or delayed.

## 7 Conclusion

In this paper, we have shown that when very little expert data is available, it is advantageous to imitate the higher-level strategy of the expert, rather than imitating the expert on action level.

To prove that the agent really imitated the strategy and not the expert actions, we considered a special setting of Imitation Learning characterized by incomplete expert trajectories. Moreover, no expert actions were available to the agent (Learning from Observations). The chosen prior for the higher-level strategy was to reach an unknown target state area. But since the expert has demonstrated how to reach it, the agent tries to reach similar states as the expert.

The presented algorithm Case-Based Inverse Reinforcement Learning (CB-IRL) builds on the idea of Temporal Coherence and Case-Based Reasoning. The algorithm trains a neural network to predict whether a state  $s_2$  can be reached from a state  $s_1$  within *windowFrame* time steps (actions). If so, the states can be considered “similar”. During inference, the agent uses this network to compare its current state  $o$  against expert states  $o_e^{(i)}$  from a Case Base. If a similar expert state  $o_e^{(j)}$  exists, the position  $j$  of this expert state in the expert trajectory serves as a (positive) reward signal for the agent. If no similar expert state exists, the agent receives a penalty. Thus, the agent is trained to reach similar states to the expert states. We tested our approach on four typical domains of Reinforcement Learning, where in every case only one tenth of an expert trajectory was available to the agent. The results show that CB-IRL was able to learn a near-optimal policy, often better than GAIfo, which had access to the full expert trajectory and was trying to imitate the expert at action level.

## References

- AD21. Saurabh Arora and Prashant Doshi. A survey of inverse reinforcement learning: Challenges, methods and progress. *Artificial Intelligence*, 297:103500, 2021.
- ADBB17. Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. A brief survey of deep reinforcement learning. *arXiv preprint arXiv:1708.05866*, 2017.
- ALUHMA08. Bryan Auslander, Stephen Lee-Urban, Chad Hogg, and Héctor Muñoz-Avila. Recognizing the enemy: Combining reinforcement learning with strategy selection using case-based reasoning. In *European Conference on Case-Based Reasoning*, pages 59–73. Springer, 2008.
- AN04. Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- BCP<sup>+</sup>16. Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- BRLdM09. Reinaldo AC Bianchi, Raquel Ros, and Ramon Lopez de Mantaras. Improving reinforcement learning by using case based heuristics. In *International Conference on Case-Based Reasoning*, pages 75–89. Springer, 2009.
- CB16. Nabanita Choudhury and Shahin Ara Begum. A survey on case-based reasoning in medicine. *International Journal of Advanced Computer Science and Applications*, 7(8):136–144, 2016.
- CLB<sup>+</sup>17. Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- DHW21. Siyu Dai, Andreas Hofmann, and Brian Williams. Automatic curricula via expert demonstrations. *arXiv preprint arXiv:2106.09159*, 2021.
- DTLS18. Debidatta DwibediR, Jonathan Tompson, Corey LynchR, and Pierre Sermanet. Self-supervised representation learning for continuous control. 2018.
- FDH<sup>+</sup>19. Carlos Florensa, Jonas Degraeve, Nicolas Heess, Jost Tobias Springenberg, and Martin Riedmiller. Self-supervised learning of image embedding for continuous control. *arXiv preprint arXiv:1901.00943*, 2019.
- FLL17. Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- GBT<sup>+</sup>15. Ross Goroshin, Joan Bruna, Jonathan Tompson, David Eigen, and Yann LeCun. Unsupervised learning of spatiotemporally coherent metrics. In *Proceedings of the IEEE international conference on computer vision*, pages 4086–4093, 2015.
- HAE<sup>+</sup>20. Lukas Hermann, Max Argus, Andreas Eitel, Artemij Amiranashvili, Wolfram Burgard, and Thomas Brox. Adaptive curriculum generation from demonstrations for sim-to-real visuomotor control. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 6498–6505. IEEE, 2020.
- HE16. Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.

- KHW<sup>+</sup>21. Joshua Knights, Ben Harwood, Daniel Ward, Anthony Vanderkop, Olivia Mackenzie-Ross, and Peyman Moghadam. Temporally coherent embeddings for self-supervised video representation learning. In *25th International Conference on Pattern Recognition (ICPR)*, pages 8914–8921. IEEE, 2021.
- KK19. Mark T Keane and Eoin M Kenny. How case-based reasoning explains neural networks: A theoretical analysis of XAI using post-hoc explanation-by-example from a survey of ANN-CBR twin-systems. In *International Conference on Case-Based Reasoning*, pages 155–171. Springer, 2019.
- LHYL19. Youngwoon Lee, Edward S Hu, Zhengyu Yang, and Joseph J Lim. To follow or not to follow: Selective imitation learning from observations. *arXiv preprint arXiv:1912.07670*, 2019.
- LLCR18. Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- LPS21. Wei Li, Florentina Paraschiv, and Georgios Sermpinis. A data-driven explainable case-based reasoning approach for financial risk detection. *Available at SSRN 3912753*, 2021.
- LSSL21. Youngwoon Lee, Andrew Szot, Shao-Hua Sun, and Joseph J Lim. Generalizable imitation learning from observation via inferring goal proximity. *Advances in Neural Information Processing Systems*, 34, 2021.
- MCW09. Hossein Mobahi, Ronan Collobert, and Jason Weston. Deep learning from temporal coherence in video. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 737–744, 2009.
- MKS<sup>+</sup>15. Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Ota20. Kei Ota. Tf2rl. <https://github.com/keiohta/tf2rl/>, 2020.
- Pom91. Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991.
- RGB11. Stéphane Ross, Geoffrey Gordon, and Drew 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.
- SLC<sup>+</sup>18. Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 1134–1141. IEEE, 2018.
- SRM<sup>+</sup>18. Nikolay Savinov, Anton Raichuk, Raphaël Marinier, Damien Vincent, Marc Pollefeys, Timothy Lillicrap, and Sylvain Gelly. Episodic curiosity through reachability. *arXiv preprint arXiv:1810.02274*, 2018.
- SXL16. Pierre Sermanet, Kelvin Xu, and Sergey Levine. Unsupervised perceptual rewards for imitation learning. *arXiv preprint arXiv:1612.06699*, 2016.
- TWS18. Faraz Torabi, Garrett Warnell, and Peter Stone. Generative adversarial imitation from observation. *arXiv preprint arXiv:1807.06158*, 2018.

- WW14. Stefan Wender and Ian Watson. Combining case-based reasoning and reinforcement learning for unit navigation in real-time strategy game ai. In *International Conference on Case-Based Reasoning*, pages 511–525. Springer, 2014.
- YMH<sup>+</sup>19. Chao Yang, Xiaojian Ma, Wenbing Huang, Fuchun Sun, Huaping Liu, Junzhou Huang, and Chuang Gan. Imitation learning from observations by minimizing inverse dynamics disagreement. *Advances in neural information processing systems*, 32, 2019.
- ZNY11. Will Y Zou, Andrew Y Ng, and Kai Yu. Unsupervised learning of visual invariance with temporal coherence. In *NIPS 2011 workshop on deep learning and unsupervised feature learning*, volume 3, 2011.

## A Appendix

The OpenAI domain Half Cheetah does not normally contain any absolute position information. To make this domain a state-reaching domain, we additionally added the x-position of the Cheetah to the otherwise 17-dimensional state space. Furthermore, the action space of this domain is originally continuous, which greatly complicates exploration. To facilitate exploration, we created a modified version called “Half Cheetah Discrete”. For this, 20 random (continuous) action vectors were sampled from the continuous action space. These 20 action vectors can be seen as basis vectors of the original continuous action space and together they now form a discrete action space (consisting of 20 possible actions). If one of the 20 discrete actions is selected, the corresponding continuous action vector is executed in the background.