

---

# Bad Habits: Policy Confounding and Out-of-Trajectory Generalization in RL

---

Miguel Suau    Matthijs T. J. Spaan    Frans A. Oliehoek  
Delft University of Technology  
{m.suaudecastro, m.t.j.spaan, f.a.oliehoek}@tudelft.nl

## Abstract

Reinforcement learning agents may sometimes develop habits that are effective only when specific policies are followed. After an initial exploration phase in which agents try out different actions, they eventually converge toward a particular policy. When this occurs, the distribution of state-action trajectories becomes narrower, and agents start experiencing the same transitions again and again. At this point, spurious correlations may arise. Agents may then pick up on these correlations and learn state representations that do not generalize beyond the agent’s trajectory distribution. In this paper, we provide a mathematical characterization of this phenomenon, which we refer to as policy confounding, and show, through a series of examples, when and how it occurs in practice.

## 1 Introduction

*This morning, I went to the kitchen for a coffee. When I arrived,  
I forgot why I was there, so I got myself a coffee—*

How often do you do something without paying close attention to your actions? Have you ever caught yourself thinking about something else while washing the dishes, making coffee, or cycling? Acting out of habit is a vital human skill as it allows us to concentrate on more important matters while carrying out routine tasks. You can commute to work while thinking about how to persuade your boss to give you a salary raise or prepare dinner while imagining your next holidays in the Alps. However, unlike in the above example, habits can also lead to undesired outcomes when we fail to recognize that the context has changed. You may hop in your car and start driving towards work even though it is a Sunday and you actually want to go to the grocery store, or you may flip the light switch when leaving a room even though the lights are already off.

Here we show how reinforcement learning (RL) agents may also suffer from this phenomenon. Agents can exploit spurious correlations (Pearl et al., 2016) between observed variables and rewards to build simple habits that require little effort to carry out. Such correlations are induced by the agent’s policy and hence can be relied upon so long as said policy is followed consistently. However, as we shall see, even minor trajectory deviations can result in catastrophic outcomes. Ideally, the agent should only pick up on correlations that are stable across policies. That is, independently of the trajectories being followed. We refer to this objective as *out-of-trajectory* (OOT) generalization.

**Contributions** This paper characterizes *policy confounding*, a term we use to name the above-described phenomenon. To do so, we introduce a mathematical framework that helps us investigate different types of state representations. Moreover, we provide a series of clarifying examples that illustrate how, as a result of policy confounding, the agent may learn representations based on spurious correlations that do not guarantee OOT generalization. Unfortunately, we do not have a complete answer for how to prevent policy confounding. However, we suggest a few off-the-shelf solutions that

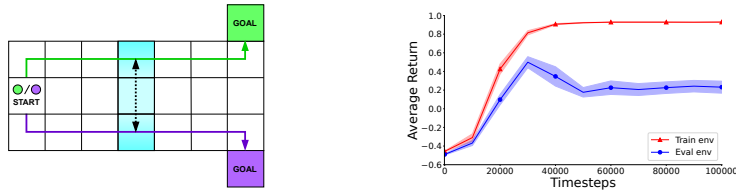


Figure 1: Left: An illustration of the Frozen T-Maze environment. Right: Learning curves when evaluated in the Frozen T-Maze environment with (blue curve) and without (red curve) ice.

may help mitigate its effects. We hope this paper will create awareness among the RL community about the risks of policy confounding and inspire further research on this topic.

## 2 Example: Frozen T-Maze

We now provide an example to illustrate the phenomenon of policy confounding and motivate the need for careful analysis. The environment shown in Figure 1 is a variant of the popular T-Maze environment (Bakker, 2001). The agent receives a binary signal, green or purple, at the start location. Then, it needs to move to the right and reach the correct goal at the end of the maze (ignore the blue cells and the black vertical arrow in the middle of the maze for now). The agent obtains a reward of  $+1$  for moving to the green (purple) goal when having received the green (purple) signal and a reward of  $-1$  otherwise. At first sight, one may think that the only way the agent can solve the task is if, at every cell along its trajectory, it can recall the initial signal. However, once the agent figures out the shortest path to each of the two goals (depicted by the green and purple arrows), the agent may safely forget the initial signal. The agent knows that whenever it is at any of the cells along the green (purple) path, it must have received the green (purple) signal. Hence, it can simply move toward the right goal on the basis of its own location. Sticking to this habit is optimal so long as the agent commits to always taking these two paths.<sup>1</sup> It is also essential that the environment’s dynamics remain the same since even the slightest change in the agent’s trajectories may erase the spurious correlation induced by the agent’s policy between the agent’s location and the correct goal. To show that this actually occurs in practice, we train agents in the original environment (train env) and evaluate them on a variant of the same (eval env), where some ice (blue) has appeared in the middle of the maze. The ice makes the agent slip from the upper cell to the bottom cell and vice versa. The plot on the right of Figure 1 shows the return averaged over 10 trials. The performance drop in the evaluation environment (blue curve) suggests that the agents’ policies do not generalize. The ice confuses the agents, who, after being pushed away from their preferred trajectories, can no longer select the right goal. More details about this experiment are provided in Section 7.

## 3 Related Work

The presence of spurious correlations in the training data is a well-studied problem in machine learning. These correlations often provide convenient shortcuts that a model can exploit to make predictions (Beery et al., 2018). However, the performance of a model that relies on them may significantly deteriorate under different data distributions (Quionero-Candela et al., 2009; Arjovsky, 2021). Langosco et al. (2022) show that RL agents may use certain environment features as proxies for choosing their actions. These features, which show only in the training environments, happen to be spuriously correlated with the agent’s objectives. In contrast, we demonstrate that, as a result of policy confounding, agents may directly take part in the formation of spurious correlations. A few prior works have already reported empirical evidence of particular forms of policy confounding, showing that in deterministic environments, agents can rely on information that correlates with the agent’s progress in an episode to determine the optimal actions. This strategy is effective because under fixed policies, features such as timers (Song et al., 2020), agent’s postures (Lan et al., 2023), or previous action sequences (Machado et al., 2018) can be directly mapped to the agent’s state. These works provide various hypotheses to justify their experimental observations. Here, we contribute an overarching theory that explains the underlying causes and mechanisms behind these results, along with a series of examples illustrating other types of policy confounding. Please refer to Appendix C for more details on related work.

<sup>1</sup>Note that the two paths highlighted in Figure 1 are not the only optimal paths. However, for the agent to be able to ignore the initial signal, it is important that the paths do not overlap.

## 4 Preliminaries

Although, as we shall see in the experiments, policy confounding can occur even when states are fully observable, in order to understand the idea, it is useful to formulate the setting as partially observable (Kaelbling et al., 1996). Moreover, since we model values and policies using (parametric) functions rather than tables, we use state variables or state factors to represent the different states of the environment (Boutilier et al., 1999).

**Definition 1** (FPOMDP). A factored partially observable Markov decision process (FPOMDP) is a tuple  $\langle S, X, A, T, R, O \rangle$ , where  $S$  is the set of state variables  $S = \{S^1, \dots, S^{|S|}\}$  defining the different states of the environment  $s = \langle s^1, \dots, s^{|S|} \rangle \in \times_i \text{dom}(S^i)$ ,  $X = \{X^1, \dots, X^{|X|}\} \subseteq S$  is the subset of state variables defining the agent’s observations  $x = \langle x^1, \dots, x^{|X|} \rangle \in \times_i \text{dom}(X^i)$ ,  $A$  is a random variable denoting the actions  $a \in \text{dom}(A)$  that are available to the agent,  $T(s_{t+1} | s_t, a_t)$  is the transition probability,  $R(s_t, a_t)$  is the immediate reward,  $O(x_t | s_t)$  is the observation probability, with  $O$  being an indicator function determining whether or not  $x_t$  is feasible given  $s_t$ .

In this setting, the agent must keep track of past actions and observations to make the right action choices (Singh et al., 1994). Policies are therefore mappings from the past action-observation history,  $h_t = \langle x_1, a_1, \dots, a_{t-1}, x_t \rangle$ , to a probability distribution over actions  $A$ ,  $\pi : H \rightarrow \Delta(A)$ , where  $H$  is the set of all possible histories of any length. We use the random variable  $\tau = \langle x_1, a_1, \dots, a_{T-1}, x_K \rangle$  to denote the agent’s trajectory in an episode, with  $K$  being the episode’s horizon. Knowing that the full history constitutes a Markov representation, we can reformulate the FPOMDP into a factored history MDP (FHMDP).

**Definition 2** (FHMPD). A factored history Markov decision process (FHMDP) is a tuple  $\langle H, \Theta, A, T_h, R_h \rangle$ , where  $H$  is the set of all possible histories of any length,  $\Theta$  denotes the set of history variables, with  $\Theta_t$  denoting the set of actions  $A$  and observation variables  $X$  in a history of length  $t$ ,  $\Theta_t = \{X_1^1, \dots, X_1^{|X|}, A_1, \dots, X_t^1, \dots, X_t^{|X|}, A_t\}$ , such that the set of histories of length  $t$  is defined as  $H_t = \times_i \text{dom}(\Theta_t^i)$ ,  $A$  is a random variable denoting the actions  $a \in \text{dom}(A)$  that are available to the agent,

$$T_h(h_{t+1} = \langle h_t, a_t, x_{t+1} \rangle | h_t, a_t) \triangleq \sum_{s_{t+1}, s_t \in S} O(x_{t+1} | s_{t+1}) T(s_{t+1} | s_t, a_t) \Pr(s_t | h_t)$$

is the history transition probability,<sup>2</sup> and

$$R_h(h_t, a_t) \triangleq \sum_{s_t \in S} R(s_t, a_t) \Pr(s_t | h_t)$$

is the history reward.

This formulation is convenient because it allows solving the POMDP using MDP methods. Yet, due to combinatorial explosion, learning a policy that conditions on the full history is generally infeasible. Fortunately, in many problems, not all the information is strictly relevant; the agent can usually find compact representations of the history, that are sufficient for solving the task (McCallum, 1995).

## 5 History representations

Factored representations are useful because they readily define relationships between (states) histories. Histories can be compared to one another by looking at the individual values the different variables take. Removing some of the variables in  $\Theta_t$  has the effect of grouping together those histories that share the same values for the remaining ones. Thus, in contrast with most of the theoretical work in RL, which treats histories (states) as independent entities, we can define history (state) abstractions at the variable level instead of doing so at the history (state) level (Li et al., 2006).

**Definition 3** (History representation). A history representation is a function  $\Phi : H_t \rightarrow \bar{H}_t$ , with  $H_t = \times_i \text{dom}(\Theta_t^i)$ ,  $\bar{H}_t = \times_i \text{dom}(\bar{\Theta}_t^i)$ , and  $\bar{\Theta}_t \subseteq \Theta_t$ .

Intuitively a history representation  $\Phi(h_t)$  is a context-specific projection of a history  $h_t \in H_t = \times_i \text{dom}(\Theta_t^i)$  onto a lower dimensional space  $\bar{H}_t = \times_i \text{dom}(\bar{\Theta}_t^i)$  defined by a subset of its variables,  $\bar{\Theta}_t \subseteq \Theta_t$ . We use  $\{h_t\}^\Phi = \{h'_t \in H_t : \Phi(h'_t) = \Phi(h_t)\}$  to denote  $h_t$ ’s equivalence class under  $\Phi$ .

<sup>2</sup>Note that we sum over  $s_{t+1}$  because multiple states may emit the same observation  $x_{t+1}$ .

## 5.1 Markov history representations

As noted in Section 4, the agent should strive for history representations with few variables. Yet, not all history representations will be sufficient to learn the optimal policy; some may exclude variables that contain useful information for the task at hand.

**Definition 4** (Markov history representation). A history representation  $\Phi(h_t)$  is said to be Markov if, for all  $h_t, h_{t+1} \in H, a_t \in A$ ,

$$R_h(h_t, a_t) = R_h(\Phi(h_t), a_t) \quad \text{and} \quad \sum_{h'_{t+1} \in \{h_{t+1}\}^\Phi} T_h(h'_{t+1} | h_t, a_t) = \Pr(\Phi(h_{t+1}) | \Phi(h_t), a_t),$$

where  $R_h(\Phi(h_t), a_t) = \{R(h'_t, a_t)\}_{h'_t \in \{h_t\}^\Phi}$  is the reward at any  $h'_t \in \{h_t\}^\Phi$ .

The above definition is equivalent to the notion of bisimulation (Dean and Givan, 1997; Givan et al., 2003) or model-irrelevance state abstraction (Li et al., 2006). Representations satisfying these conditions are guaranteed to be equivalent to the original representation. That is, for any given policy and initial history, the expected return (i.e., cumulative reward; Sutton and Barto, 2018) is the same when conditioning on the full history or on the Markov history representation. Note that a history representation  $\Phi$  such that  $\Phi(h_t) = h_t$ , for all  $h_t \in H$ , is, in itself, Markov.

**Definition 5** (Minimal history representation). A history representation  $\Phi^* : H_t \rightarrow \bar{H}_t^*$  with  $\bar{H}_t^* = \times_i \text{dom}(\bar{\Theta}_t^{*i})$  is said to be *minimal*, if all other history representations  $\Phi : H_t \rightarrow \bar{H}_t$  with  $\bar{H}_t = \times_i \text{dom}(\bar{\Theta}_t^i)$  and  $|\bar{\Theta}_t| \subset |\bar{\Theta}_t^*|$ , for at least one  $h_t \in H$ , are not Markov.

In other words,  $\Phi_t^*(h_t)$  is *minimal* when none of the remaining variables can be removed while the representation remains Markov. Hence, we say that a minimal history representation  $\Phi_t^*(h_t)$  is a sufficient statistic of the full history.

**Definition 6** (Superfluous variable). Let  $\{\bar{\Theta}_t^*\}_{\cup \Phi^*}$  be the union of variables in all possible minimal history representations. A variable  $\Theta_t^i \in \Theta_t$  is said to be superfluous, if  $\Theta_t^i \notin \{\bar{\Theta}_t^*\}_{\cup \Phi^*}$ .

## 5.2 $\pi$ -Markov history representations

Considering that the agent’s policy will rarely visit all possible histories, the notion of Markov history representation seems excessively strict. We now define a relaxed version that guarantees the representation to be Markov when a specific policy  $\pi$  is followed.

**Definition 7** ( $\pi$ -Markov history representation). A history representation  $\Phi^\pi(h_t)$  is said to be  $\pi$ -Markov if, for all  $h_t, h_{t+1} \in H^\pi, a_t \in \text{supp}(\pi(\cdot | h_t))$ ,

$$R_h(h_t, a_t) = R_h^\pi(\Phi^\pi(h_t), a_t) \quad \text{and} \quad \sum_{h'_{t+1} \in \{h_{t+1}\}^\Phi} T_h(h'_{t+1} | h_t, a_t) = \Pr^\pi(\Phi^\pi(h_{t+1}) | \Phi^\pi(h_t), a_t),$$

where  $H^\pi \subseteq H$  denotes the histories visited under  $\pi$ ,  $R_h^\pi(\Phi^\pi(h_t), a_t) = \{R_h(h'_t, a_t)\}_{h'_t \in \{h_t\}^\Phi}$ ,  $\{h_t\}^\Phi_\pi = \{h'_t \in H_t^\pi : \Phi^\pi(h'_t) = \Phi^\pi(h_t)\}$ , and  $\Pr^\pi$  is probability under  $\pi$ .

**Definition 8** ( $\pi$ -minimal history representation). A history representation  $\Phi^{\pi*} : H_t^\pi \rightarrow \bar{H}_t^{\pi*}$  with  $\bar{H}_t^{\pi*} = \times_i \text{dom}(\bar{\Theta}_t^{\pi*i})$  is said to be  $\pi$ -*minimal*, if all other history representations  $\Phi : H_t^\pi \rightarrow \bar{H}_t^\pi$  with  $\bar{H}_t^\pi = \times_i \text{dom}(\bar{\Theta}_t^i)$  and  $|\bar{\Theta}_t| \subset |\bar{\Theta}_t^{\pi*}|$ , for at least one  $h_t \in H^\pi$ , are not  $\pi$ -Markov.

## 6 Policy Confounding

We are now ready to describe how and when policy confounding occurs, as well as why we should care, and how we should go about preventing it. The proofs for all theoretical results are deferred to Appendix A.

Policy confounding arises naturally as the agent improves its policy. Normally, at the beginning of training, the agent takes exploratory actions to determine which ones yield high rewards. It is only after the agent has committed to a particular policy that we start seeing how some of the variables in its history become irrelevant for predicting future states and rewards. The agent may then choose to ignore these variables and exclude them from its representation if keeping them takes extra ‘effort’.

The next result demonstrates that a  $\pi$ -Markov history representation  $\Phi^\pi$  requires at most the same variables, and in some cases fewer, than a minimal history representation  $\Phi^*$ , while still satisfying the Markov conditions for those histories visited under  $\pi$ ,  $h_t \in H^\pi$ .

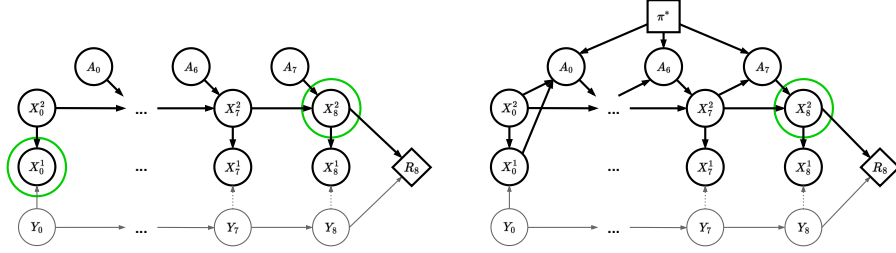


Figure 2: Two DBNs representing the dynamics of the Frozen T-Maze environment, when actions are sampled at random (left), and when they are determined by the optimal policy (right).

**Proposition 1.** Let  $\Phi^*$  be the set of all possible minimal history representations, where every  $\Phi^* \in \Phi^*$  is defined as  $\Phi^* : H_t \rightarrow \bar{H}_t^*$  with  $\bar{H}_t^* = \times_i \text{dom}(\Theta_t^{*i})$ . For all  $\pi$  and all  $\Phi^* \in \Phi^*$ , there exists a  $\pi$ -Markov history representation  $\Phi^\pi : H_t^\pi \rightarrow \bar{H}_t^\pi$  with  $\bar{H}_t^\pi = \times_i \text{dom}(\bar{\Theta}_t^{\pi i})$  such that for all  $h_t \in H_t^\pi$ ,  $\bar{\Theta}_t^\pi \subseteq \bar{\Theta}_t^*$ . Moreover, there exist cases for which  $\bar{\Theta}_t^\pi$  is a proper subset,  $\bar{\Theta}_t^\pi \neq \bar{\Theta}_t^*$ .

Although the result above seems intuitive, its truth may appear incidental. While it is clear that  $\Phi^\pi$  will never require more variables than the corresponding minimal history representation  $\Phi^*$ , whether or not  $\Phi^\pi$  will require fewer, seems just an arbitrary consequence of the policy being followed. Moreover, since the variables in  $\bar{\Theta}_t^*$  are all strictly relevant for predicting transitions and rewards, one may think that a policy  $\pi$  inducing representations such that  $\bar{\Theta}_t^\pi \subset \bar{\Theta}_t^*$  can never be optimal. However, as shown by the following example, it turns out that the histories visited by a particular policy, especially if it is the optimal policy, tend to contain a lot of redundant information. This is particularly true in environments where future observations are heavily influenced by past actions and observations. In such cases, the current observation often reveals a lot about the agent’s trajectory.

**Example 1. (Frozen T-Maze)** Let us consider the Frozen T-Maze again (Section 2). Figure 3 shows a dynamic Bayesian network (DBN; Dean and Kanazawa, 1989; Murphy, 2002) describing the dynamics of the environment. Observation variables are denoted by  $X$ , while hidden variables are denoted by  $Y$ . The nodes labeled as  $X^2$  represent the agent’s location from  $t = 0$  to  $t = 8$ . All intermediate nodes between  $t = 0$  and  $t = 7$  are omitted for simplicity. The nodes labeled as  $Y$  indicate whether the goal is to go to the green or the purple cell (see Figure 1). Note that  $Y$  always takes the same value at all timesteps within an episode (either green or purple). The information in  $Y$  is hidden and only passed to the agent at the start location through the node  $X_0^1$ . On the one hand, if actions are not specified by any particular policy, but simply sampled at random (left diagram), to determine the reward  $R_8$  at  $t = 8$ , one needs to know the signal  $X_0^1$  received at  $t = 0$  and the agent’s current location  $X_8^2$ . These are highlighted by the green circles in the left DBN. This is because the actions  $\langle A_0, \dots, A_7 \rangle$  appear as exogenous variables and can take any possible value. Hence, the reward could be either  $-0.1$ , (per timestep penalty),  $-1$  (wrong goal), or  $+1$  (correct goal) depending on the actual values of  $X_1^1$  and  $X_8^2$ . On the other hand, when actions are sampled from the optimal policy  $\pi^*$  (right DBN), knowing  $X_8^2$  (green circle) is sufficient to determine  $r_8$ . In this second case,  $\pi^*$  makes the action  $A_0$ , and thus all future agent locations, dependent on the initial signal  $X_0^1$ . This occurs because, under the optimal policy (green and purple paths in Figure 1), the agent always takes the action ‘move up’ when receiving the green signal or ‘move down’ when receiving the purple signal, and then follows the shortest path towards each of the goals. As such, we have that, from  $t = 1$  onward,  $\Phi^{\pi^*}(h_t) = X_t^2$  is a  $\pi$ -Markov history representation since it constitutes a sufficient statistic of the history  $H_t$  under  $\pi^*$ . Finally, note that, for the same reason, from  $t = 1$ , actions may also condition only on  $X^2$ .

The phenomenon highlighted by the previous example is the result of a spurious correlation induced by the optimal policy between the agent’s locations  $\langle X_0^2, \dots, X_8^2 \rangle$  and the reward  $R_8$ . Generally speaking, this occurs because policies act as confounders, opening backdoor paths between future histories  $\Theta_{t+1}$  and the variables in the current history  $\Theta_t$  (Pearl, 2000). This is shown by the DBN depicted in Figure 3, where we see that the policy influences both the current history variables and also future history variables, hence potentially affecting their conditional relationships. For instance, in the above example,  $R^{\pi^*}(X_8^2 = \text{‘green goal’}) = +1$  when following  $\pi^*$ , while for an arbitrary  $\pi$ ,  $R(X_8^2 = \text{‘green goal’}) = \pm 1$ .

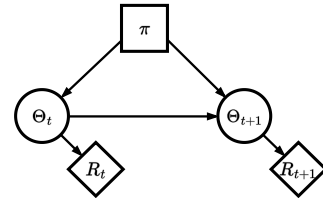


Figure 3: A DBN illustrating the phenomenon of policy confounding. The policy opens backdoor path that can affect conditional relations between the variables in  $h_t$  and  $h_{t+1}$

**Definition 9** (Policy Confounding). A history representation  $\Phi : H_t \rightarrow \bar{H}_t$  is said to be confounded by a policy  $\pi$  if, for some  $h_t, h_{t+1} \in H, a_t \in A$ ,

$$R^\pi(\Phi(h_t), a_t) \neq R^\pi(\text{do}(\Phi(h_t)), a_t) \quad \text{or} \quad \Pr^\pi(\Phi(h_{t+1}) \mid \Phi(h_t), a_t) \neq \Pr^\pi(\Phi(h_{t+1}) \mid \text{do}(\Phi(h_t)), a_t)$$

The operator  $\text{do}(\cdot)$  is known as the do-operator, and it is used to represent physical interventions in a system (Pearl, 2000). These interventions are meant to distinguish cause-effect relations from mere statistical associations. In our case,  $\text{do}(\Phi(h_t))$  means setting the variables forming the history representation  $\Phi(h_t)$  to a particular value and considering all possible histories in the equivalence class,  $h'_t \in \{h_t\}^\Phi$ . That is, independently of what policy is being followed.

It is easy to show that the underlying reason why a  $\pi$ -Markov history representation may require fewer variables than the minimal history representation (as in Example 1) is indeed policy confounding.

**Theorem 1.** *Let  $\Phi^* : H_t \rightarrow \bar{H}_t^*$  with  $\bar{H}_t^* = \times_i \text{dom}(\bar{\Theta}_t^{*i})$  be a minimal history representation. If, for some  $\pi$ , there is a  $\pi$ -Markov history representation  $\Phi^\pi : H_t^\pi \rightarrow \bar{H}_t^\pi$  with  $\bar{H}_t^\pi = \times_i \text{dom}(\bar{\Theta}_t^{\pi i})$ , such that  $\bar{\Theta}_t^\pi \subset \bar{\Theta}_t^*$  for some  $h_t \in H$ , then  $\Phi^\pi$  is confounded by policy  $\pi$ .*

Finally, to conclude this section, we demonstrate that even though, in Example 1, the variables included in the  $\pi$ -minimal history representation are a subset of the variables in the minimal history representation,  $\bar{\Theta}_t^{\pi*} \subset \bar{\Theta}_t^*$ , this is not always the case, as  $\bar{\Theta}_t^{\pi*}$  may contain superfluous variables (Definition 6). An example illustrating this situation is provided in Appendix B (Example 4).

**Proposition 2.** *Let  $\{\bar{\Theta}_t^*\}_{\cup \Phi^*}$  be the union of variables in all possible minimal history representations. There exist cases where, for some  $\pi$ , there is a  $\pi$ -minimal history representation  $\Phi^{\pi*} : H_t^\pi \rightarrow \bar{H}_t^{\pi*}$  with  $\bar{H}_t^{\pi*} = \times_i \text{dom}(\bar{\Theta}_t^{\pi*i})$  such that  $\bar{\Theta}_t^{\pi*} \setminus \{\bar{\Theta}_t^*\}_{\cup \Phi^*} \neq \emptyset$ .*

## 6.1 Why should we care about policy confounding?

Leveraging spurious correlations to develop simple habits can be advantageous when resources such as memory, computing power, or data are limited. Agents can disregard and exclude from their representation those variables that are redundant under their policies. However, the challenge is that some of these variables may be crucial to ensure that the agent behaves correctly when the context changes. In the Frozen T-Maze example from Section 2, we observed how the agent could no longer find the correct goal when the ice pushed it away from the optimal trajectory. This is a specific case of a well-researched issue known as out-of-distribution (OOD) generalization (Quionero-Candela et al., 2009; Arjovsky, 2021). We refer to it as *out-of-trajectory* (OOT) generalization to highlight that the problem arises due to repeatedly sampling from the same policy and thus following the same trajectories. In contrast to previous works (Kirk et al., 2023) that address generalization to environments that differ from the training environment, our objective here is to generalize to trajectories the agent never (or only rarely) takes.<sup>3</sup>

Ideally, the agent should aim to learn representations that enable it to predict future rewards and transitions even when experiencing slight variations in its trajectory. Based on Definition 4, we know that, in general, only a Markov history representation satisfies these requirements. However, computing such representations is typically intractable (Ferns et al., 2006), and thus most standard RL methods usually learn representations by maximizing an objective function that depends on the distribution of trajectories  $P^b(\tau)$  visited under a behavior policy  $b$  (e.g., expected return,  $\mathbb{E}_{\tau \sim P^b(\tau)} [G(\tau)]$ ; Sutton and Barto, 2018). The problem is that  $b$  may favor certain trajectories over others, which may lead to the exploitation of spurious correlations in the learned representation.

## 6.2 When should we worry about OOT generalization in practice?

The previous section highlighted the generalization failures of representations that depend on spurious correlations. Now, let us delve into the circumstances in which policy confounding is most prone to cause problems.

<sup>3</sup>Note that in the Frozen T-Maze environment, the ice does change the environment dynamics. However, its purpose is to compel the agent to take trajectories different from the optimal ones. The way we implemented it, the effect of the ice would be equivalent to forcing the agent to move down twice when in the top cell or move up twice when in the bottom cell. These trajectories are feasible in the original environment.

**Function approximation** Function approximation has enabled traditional RL methods to scale to high-dimensional problems with long-term memory dependencies, where storing values in lookup tables is infeasible. Using parametric functions (e.g., neural networks) to model policies and value functions, agents can learn abstractions by grouping together histories if these yield the same transitions and rewards. As mentioned before, abstractions occur naturally when histories are represented by a set of variables since the functions simply need to ignore some of these variables. However, this also implies that value functions and policies are exposed to spurious correlations. If a particular variable becomes irrelevant due to policy confounding, the function may learn to ignore it and remove it from its representation (Example 1). This is in contrast to tabular representations, where, every history takes a separate entry, and even though there exist algorithms that perform history (state) abstractions in tabular settings (Andre and Russell, 2002; Givan et al., 2003), these abstractions are normally formed offline before learning (computing) the policy, hence avoiding the risk of policy confounding.

**Narrow trajectory distributions** In practice, agents are less prone to policy confounding when the trajectory distribution  $P^b(\tau)$  is broad (i.e., when  $b$  encompasses a wide set of trajectories) than when it is narrow. This is because the spurious correlations present in certain trajectories are less likely to have an effect on the learned representations. On-policy methods (e.g., SARSA, Actor-Critic; Sutton and Barto, 2018) are particularly troublesome for this reason since the same policy being updated must also be used to collect the samples. Yet, even when the trajectory distribution is narrow, there is no reason why the agent should pick up on spurious correlations while its policy is still being updated. Only when the agent commits to a particular policy should we start worrying about policy confounding. At this point, lots of the same trajectories are being used for training, and the agent may ‘forget’ (French, 1999) that, even though certain variables may no longer be needed to represent the current policy, they were important under previous policies. This generally occurs at the end of training when the agent has converged to a particular policy. However, if policy confounding occurs earlier during training, it may prevent the agent from further improving its policy (Nikishin et al., 2022; please refer to Appendix C for more details).

### 6.3 What can we do to improve OOT generalization?

As mentioned in the introduction, we do not have a complete answer to the problem of policy confounding. Yet, here we offer a few off-the-shelf solutions that, while perhaps limited in scope, can help mitigate the problem in some situations. These solutions revolve around the idea of broadening the distribution of trajectories so as to dilute the spurious correlations introduced by certain policies.

**Off-policy methods** We already explained in Section 6.2 that on-policy methods are particularly prone to policy confounding since they are restricted to using samples coming from the same policy. A rather obvious solution is to instead use off-policy methods, which allow using data generated from previous policies. Because the samples belong to a mixture of policies it is less likely that the model will pick up the spurious correlations present on specific trajectories. However, as we shall see in the experiments, this alternative works only when replay buffers are large enough. This is because standard replay buffers are implemented as queues, and hence the first experiences coming in are the first being removed. This implies that a replay buffer that is too small will contain samples coming from few and very similar policies. Since there is a limit on how large replay buffers are allowed to be, future research could explore other, more sophisticated, ways of deciding what samples to store and which ones to remove (Schaul et al., 2016).

**Exploration and domain randomization** When allowed, exploration may mitigate the effects of policy confounding and prevent agents from overfitting their preferred trajectories. Exploration strategies have already been used for the purpose of generalization; to guarantee robustness to perturbations in the environment dynamics (Eysenbach and Levine, 2022), or to boost generalization to unseen environments (Jiang et al., 2022). The goal for us is to remove, to the extent possible, the spurious correlations introduced by the current policy. Unfortunately, though, exploration is not always without cost. Safety-critical applications require the agent to stay within certain boundaries (Altman, 1999; García and Fernández, 2015). When training on a simulator, an alternative to exploration is domain randomization (Tobin et al., 2017; Peng et al., 2018; Machado et al., 2018). The empirical results reported in the next section suggest that agents become less susceptible to policy confounding when adding enough stochasticity to the environment or to the policy. Yet, there is a



Figure 4: Illustrations of the Key2Door (left) and Diversion (right) environments.

limit on how much noise can be added to the environment or the policy without altering the optimal policy ( Sutton and Barto, 2018, Example 6.6: Cliff Walking).

## 7 Experiments

The goal of the experiments is to: (1) demonstrate that the phenomenon of policy confounding described by the theory does occur in practice, (2) uncover the circumstances under which agents are most likely to suffer the effects of policy confounding and fail to generalize, and (3) evaluate how effective the strategies proposed in the previous section are in mitigating these effects.

### 7.1 Experimental setup

Agents are trained with an off-policy method, DQN (Mnih et al., 2015) and an on-policy method, PPO (Schulman et al., 2017). To be able to analyze the learned representations more easily, we represent policies and value functions as feedforward neural networks and use a stack of past observations as input in the environments that require memory. We report the mean return as a function of the number of training steps. Training is interleaved with periodic evaluations on the original environments and variants thereof used for validation. The results are averaged over 10 random seeds. Please refer to Appendix F for more details about the experimental setup.

### 7.2 Environments

We ran our experiments on three grid-world environments: the **Frozen T-Maze** from Section 2, and the below described **Key2Door**, and **Diversion** environments. We use these as pedagogical examples to clarify the ideas introduced by the theory. Nonetheless, in Appendix C, we refer to previous works showing evidence of particular forms of policy confounding in high dimensional domains.

**Example 2. Key2Door.** Here, the agent needs to collect a key placed at the beginning of the corridor in Figure 4 (left) and then open the door at the end. The observations do not show whether the key has already been collected. Thus, to solve the task in the minimum number of steps, the agent must remember that it already got the key when going to the door. Yet, since during training, the agent always starts the episode at the first cell from the left, when moving towards the door, the agent can forget about the key once it has reached the third cell. As in the Frozen T-Maze example, the agent can build the habit of using its own location to tell whether it has or has not got the key yet. This, can only occur when the agent consistently follows the optimal policy, depicted by the purple arrow. Otherwise, if the agent moves randomly through the corridor, it is impossible to tell whether the key has or has not been collected. In contrast, in the evaluation environment, the agent always starts at the second to last cell, this confuses the agent, which is used to already having the key by the time it reaches said cell. A DBN describing the dynamics of the environment is provided in Appendix D.

**Example 3. Diversion.** Here, the agent must move from the start state to the goal state in Figure 4 (right). The observations are length-8 binary vectors. The first 7 elements indicate the column where the agent is located. The last element indicates the row. This environment aims to show that policy confounding can occur not only when the environment is partially observable, as was the case in the previous examples, but also in fully observable scenarios. After the agent learns the optimal trajectory depicted by the green arrow, it can disregard the last element in the observation vector. This is because, if the agent does not deviate, the bottom row is never visited. Rather than forgetting past information, the agent ignores the last element in the current observation vector for being irrelevant when following the optimal trajectory. We train the agent in the original environment and evaluate it in a version with a yellow diversion sign in the middle of the maze that forces the agent to move to the bottom row. A DBN describing the dynamics of the environment is provided in Appendix D.



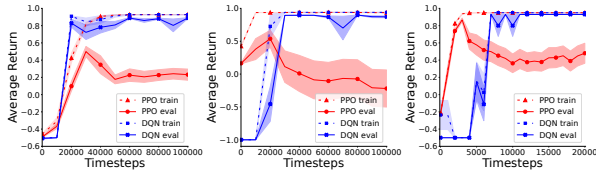


Figure 5: DQN vs. PPO in the train and evaluation variants of Frozen T-Maze (left), Key2Door (middle), and Diversion (right).

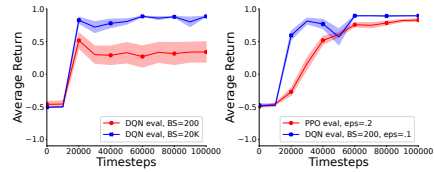


Figure 6: Frozen T-Maze. Left: DQN small vs. large buffer sizes. Right: PPO and DQN when adding stochasticity.

### 7.3 Results

**On-policy vs. off-policy** The results in Figure 5 reveal the same pattern in all three environments. PPO fails to generalize outside the agent’s preferred trajectories. After an initial phase where the average returns on the training and evaluation environments increase (‘PPO train’ and ‘PPO eval’), the return on the evaluation environments (‘PPO eval’) starts decreasing when the agent commits to a particular trajectory, as a result of policy confounding. In contrast, since the training samples come from a mixture of policies, DQN performs optimally in both variants of the environments (‘DQN train’ and ‘DQN eval’) long after converging to the optimal policy.<sup>4</sup> A visualization of the history representations learned with PPO, showing that the policy does ignore variables that are necessary for generalization, is provided in Appendix E.1.

**Large vs. small replay buffers** We mentioned in Section 6.3 that the effectiveness of off-policy methods against policy confounding depends on the size of the replay buffer. The results in Figure 6 (left) confirm this claim. The plot shows the performance of DQN in the Frozen T-Maze environment when the size of the replay buffer contains 100K experiences and when it only contains the last 10K experiences. We see that in the second case, the agents performance in the evaluation environment decreases (red curve left plot). This is because, after the initial exploration phase, the distribution of trajectories becomes too narrow, and the spurious correlations induced by the latest policies dominate the replay buffer. Similar results for the other two environments are provided in Appendix E.2.

**Exploration and domain randomization** The last experiment shows that if sufficient exploration is allowed, DQN may still generalize to different trajectories, even when using small replay buffers (blue curve right plot Figure 6). In the original configuration, the exploration rate  $\epsilon$  for DQN starts at  $\epsilon = 1$  and decays linearly to  $\epsilon = 0.0$  after 20K steps. For this experiment, we set the final exploration rate  $\epsilon = 0.1$ . In contrast, since exploration in PPO is normally controlled by the entropy bonus, which makes it hard to ensure fixed exploration rates, we add noise to the environment instead. The red curve in Figure 6 (right) shows that when we train in an environment where the agent’s actions are overridden by a random action with 20% probability, the performance of PPO in the evaluation environment does not degrade after the agent has converged to the optimal policy. This suggests that the added noise prevents the samples containing spurious correlations from dominating the training batches. However, it may also happen that random noise is not sufficient to remove the spurious correlations. As shown in Figure 13 (Appendix E.2), in the Key2Door environment, neither forcing the agent to take random actions 20% of the time nor setting  $\epsilon = 0.1$ , solves the OOT generalization problem. Similar results for Diversion are provided in Appendix E.2.

## 8 Conclusion

This paper described the phenomenon of policy confounding. We showed both theoretically and empirically how as a result of following certain trajectories, agents may pick up on spurious correlations, and build habits that are not robust to trajectory deviations. We also uncovered the circumstances under which policy confounding is most likely to occur in practice and suggested a few ad hoc solutions that may mitigate its effects. We conceive this paper as a stepping stone to explore more sophisticated solutions. An interesting avenue for future research is the integration of tools from the field of causal inference (Pearl et al., 2016; Peters et al., 2017) to aid the agent in forming history representations that are grounded on causal relationships rather than mere statistical associations (Lu et al., 2018; Zhang et al., 2020; Sontakke et al., 2021; Saengkyongam et al., 2023).

<sup>4</sup>The small gap between ‘DQN train’ and ‘DQN eval’ is due to the  $-0.1$  penalty per timestep. In all three environments, the shortest path is longer in the evaluation environment than in the training environment.

## Acknowledgements

This project received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 758824 —INFLUENCE)



## References

- Altman, E. (1999). *Constrained Markov decision processes*, volume 7. CRC press.
- Andre, D. and Russell, S. J. (2002). State abstraction for programmable reinforcement learning agents. In *Proceedings of the Eighteenth National Conference on Artificial Intelligence*, pages 119–125.
- Arjovsky, M. (2021). Out of distribution generalization in machine learning. *arXiv preprint arXiv:2103.02667*.
- Bakker, B. (2001). Reinforcement learning with long short-term memory. *Advances in neural information processing systems*, 14.
- Beery, S., Van Horn, G., and Perona, P. (2018). Recognition in terra incognita. In *Proceedings of the European conference on computer vision (ECCV)*, pages 456–473.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The Arcade Learning Environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.
- Boutillier, C., Dean, T., and Hanks, S. (1999). Decision-theoretic planning: Structural assumptions and computational leverage. *Journal of Artificial Intelligence Research*, 11:1–94.
- Castro, P. S. (2020). Scalable methods for computing state similarity in deterministic markov decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Dean, T. and Givan, R. (1997). Model minimization in Markov decision processes. In *Proc. of the National Conference on Artificial Intelligence*, pages 106–111.
- Dean, T. and Kanazawa, K. (1989). A model for reasoning about persistence and causation. *Computational intelligence*, 5(2):142–150.
- Eysenbach, B. and Levine, S. (2022). Maximum entropy RL (provably) solves some robust RL problems. In *International Conference on Learning Representations*.
- Ferns, N., Castro, P. S., Precup, D., and Panangaden, P. (2006). Methods for computing state similarity in markov decision processes. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence, UAI’06*, page 174–181.
- French, R. M. (1999). Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135.
- García, J. and Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1):1437–1480.
- Givan, R., Dean, T., and Greig, M. (2003). Equivalence notions and model minimization in Markov decision processes. *Artificial Intelligence*, 14(1–2):163–223.
- Higgins, I., Pal, A., Rusu, A., Matthey, L., Burgess, C., Pritzel, A., Botvinick, M., Blundell, C., and Lerchner, A. (2017). Darla: Improving zero-shot transfer in reinforcement learning. In *International Conference on Machine Learning*, pages 1480–1490. PMLR.
- Jiang, Y., Kolter, J. Z., and Raileanu, R. (2022). Uncertainty-driven exploration for generalization in reinforcement learning. In *Deep Reinforcement Learning Workshop NeurIPS 2022*.
- Kaelbling, L. P., Littman, M., and Moore, A. (1996). Reinforcement learning: A survey. *Journal of AI Research*, 4:237–285.

- Kirk, R., Zhang, A., Grefenstette, E., and Rocktäschel, T. (2023). A survey of zero-shot generalisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76:201–264.
- Lan, L.-C., Zhang, H., and Hsieh, C.-J. (2023). Can agents run relay race with strangers? generalization of RL to out-of-distribution trajectories. In *The Eleventh International Conference on Learning Representations*.
- Langosco, L., Koch, J., Sharkey, L. D., Pfau, J., and Krueger, D. (2022). Goal misgeneralization in deep reinforcement learning. In *International Conference on Machine Learning*, pages 12004–12019. PMLR.
- Lazaric, A. (2012). Transfer in reinforcement learning: a framework and a survey. *Reinforcement Learning: State-of-the-Art*, pages 143–173.
- Li, L., Walsh, T. J., and Littman, M. L. (2006). Towards a unified theory of state abstraction for MDPs. In *International Symposium on Artificial Intelligence and Mathematics (ISAIM 2006)*.
- Lu, C., Schölkopf, B., and Hernández-Lobato, J. M. (2018). Deconfounding reinforcement learning in observational settings. *arXiv preprint arXiv:1812.10576*.
- Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M. (2018). Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562.
- Mandlekar, A., Zhu, Y., Garg, A., Fei-Fei, L., and Savarese, S. (2017). Adversarially robust policy learning: Active construction of physically-plausible perturbations. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3932–3939. IEEE.
- McCallum, A. K. (1995). *Reinforcement Learning with Selective Perception and Hidden State*. PhD thesis, University of Rochester.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529.
- Muller-Brockhausen, M., Preuss, M., and Plaat, A. (2021). Procedural content generation: Better benchmarks for transfer reinforcement learning. In *2021 IEEE Conference on games (CoG)*, pages 01–08. IEEE.
- Murphy, K. P. (2002). *Dynamic Bayesian Networks: Representation, Inference and Learning*. PhD thesis, UC Berkeley, Computer Science Division.
- Nikishin, E., Schwarzer, M., D’Oro, P., Bacon, P.-L., and Courville, A. (2022). The primacy bias in deep reinforcement learning. In *International Conference on Machine Learning*, pages 16828–16847. PMLR.
- Ornia, D. J., Romao, L., Hammond, L., Mazo Jr, M., and Abate, A. (2022). Observational robustness and invariances in reinforcement learning via lexicographic objectives. *arXiv preprint arXiv:2209.15320*.
- Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press.
- Pearl, J., Glymour, M., and Jewell, N. P. (2016). Causal inference in statistics: A primer. 2016. *Internet resource*.
- Peng, X. B., Andrychowicz, M., Zaremba, W., and Abbeel, P. (2018). Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3803–3810. IEEE.
- Peters, J., Janzing, D., and Schölkopf, B. (2017). *Elements of causal inference: foundations and learning algorithms*. The MIT Press.
- Quionero-Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. (2009). *Dataset shift in machine learning*. The MIT Press.

- Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N. (2021). Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8.
- Saengkyongam, S., Thams, N., Peters, J., and Pfister, N. (2023). Invariant policy learning: A causal perspective. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Schaul, T., Quan, J., Antonoglou, I., and Silver, D. (2016). Prioritized experience replay. In *International Conference on Learning Representations*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Singh, S. P., Jaakkola, T., and Jordan, M. I. (1994). Learning without state-estimation in partially observable Markovian decision processes. In *Proc. of the International Conference on Machine Learning*, pages 284–292.
- Song, X., Jiang, Y., Tu, S., Du, Y., and Neyshabur, B. (2020). Observational overfitting in reinforcement learning. In *International Conference on Learning Representations*.
- Sontakke, S. A., Mehrjou, A., Itti, L., and Schölkopf, B. (2021). Causal curiosity: RL agents discovering self-supervised experiments for causal representation learning. In *International conference on machine learning*, pages 9848–9858. PMLR.
- Stone, A., Ramirez, O., Konolige, K., and Jonschkowski, R. (2021). The distracting control suite—a challenging benchmark for reinforcement learning from pixels. *arXiv preprint arXiv:2101.02722*.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Taylor, M. E. and Stone, P. (2009). Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(7).
- Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., and Abbeel, P. (2017). Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE.
- Zhang, A., Ballas, N., and Pineau, J. (2018a). A dissection of overfitting and generalization in continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*.
- Zhang, A., Lyle, C., Sodhani, S., Filos, A., Kwiatkowska, M., Pineau, J., Gal, Y., and Precup, D. (2020). Invariant causal prediction for block mdps. In *International Conference on Machine Learning*, pages 11214–11224. PMLR.
- Zhang, C., Vinyals, O., Munos, R., and Bengio, S. (2018b). A study on overfitting in deep reinforcement learning. *arXiv preprint arXiv:1804.06893*.
- Zhao, W., Queralta, J. P., and Westerlund, T. (2020). Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE symposium series on computational intelligence (SSCI)*, pages 737–744. IEEE.

## A Proofs

**Lemma 1.** Let  $\Phi^{\pi_1^*}$  be the set of all possible  $\pi$ -minimal history representations under  $\pi_1$ , where every  $\Phi^{\pi_1^*} \in \Phi^{\pi_1^*}$  is defined as  $\Phi^{\pi_1^*} : H_t^{\pi_1} \rightarrow \bar{H}_t^{\pi_1^*}$  and  $\bar{H}_t^{\pi_1^*} = \times_i \text{dom}(\bar{\Theta}_t^{\pi_1^* i})$ , and let  $\pi_2$  be a second policy such that for all  $h_t \in H_t^{\pi_1} \cap H_t^{\pi_2}$ ,

$$\text{supp}(\pi_2(\cdot | h_t)) \subseteq \text{supp}(\pi_1(\cdot | h_t)).$$

For all  $\Phi^{\pi_1^*} \in \Phi^{\pi_1^*}$ , there exists a  $\pi$ -Markov history representation under policy  $\pi_2$ ,  $\Phi^{\pi_2} : H_t^{\pi_2} \rightarrow \bar{H}_t^{\pi_2}$  with  $\bar{H}_t^{\pi_2} = \times_i \text{dom}(\bar{\Theta}_t^{\pi_2 i})$ , such that  $\bar{\Theta}_t^{\pi_2} \subseteq \bar{\Theta}_t^{\pi_1^*}$  for all  $h_t \in H_t^{\pi_1} \cap H_t^{\pi_2}$ . Moreover, there exist cases where  $\bar{\Theta}_t^{\pi_2} \neq \bar{\Theta}_t^{\pi_1^*}$ .

*Proof.* First, it is easy to show that

$$\forall h_t \in H, \text{supp}(\pi_2(\cdot | h_t)) \subseteq \text{supp}(\pi_1(\cdot | h_t)) \iff H^{\pi_2} \subseteq H^{\pi_1},$$

and

$$\forall h_t \in H, \text{supp}(\pi_2(\cdot | h_t)) = \text{supp}(\pi_1(\cdot | h_t)) \iff H^{\pi_2} = H^{\pi_1}.$$

In particular,  $H^{\pi_2} \subset H^{\pi_1}$  if there is at least one history  $h'_t \in H^{\pi_1} \cap H^{\pi_2}$  such that

$$\text{supp}(\pi_2(\cdot | h'_t)) \subset \text{supp}(\pi_1(\cdot | h'_t))$$

while

$$\text{supp}(\pi_2(\cdot | h_t)) = \text{supp}(\pi_1(\cdot | h_t))$$

for all other  $h_t \in H^{\pi_1} \cap H^{\pi_2}$ .

In such cases, we know that there is at least one action  $a'$  for which  $\pi_2(a'_t | h'_t) = 0$  but  $\pi_1(a'_t | h'_t) \neq 0$ . Hence, since  $h'_{t+1} = \langle h'_t, a'_t \rangle \notin H_2^{\pi_2}$  but  $h'_{t+1} \in H_1^{\pi_1}$ , we have that  $H^{\pi_2} \subset H^{\pi_1}$ . Note that, in some cases, having  $\pi_2(a'_t | h'_t) = 0$  may not only remove  $h'_{t+1}$  from  $H_2^{\pi_2}$  but possibly also subsequent histories that can only be reached from  $h'_{t+1}$ .

Further, since  $H^{\pi_2} \subset H^{\pi_1}$ , we know that, for every  $\Phi^{\pi_1^*} \in \Phi^{\pi_1^*}$ , there must be a  $\Phi^{\pi_2^*}$  that requires, at most, the same number of variables,  $\bar{\Theta}_t^{\pi_2^*} \subseteq \bar{\Theta}_t^{\pi_1^*}$  and, in some cases, fewer,  $\bar{\Theta}_t^{\pi_1^*} \neq \bar{\Theta}_t^{\pi_2^*}$  (e.g., Frozen T-Maze example). □

**Proposition 1.** Let  $\Phi^*$  be the set of all possible minimal history representations, where every  $\Phi^* \in \Phi^*$  is defined as  $\Phi^* : H_t \rightarrow \bar{H}_t^*$  with  $\bar{H}_t^* = \times_i \text{dom}(\bar{\Theta}_t^{*i})$ . For all  $\pi$  and all  $\Phi^* \in \Phi^*$ , there exists a  $\pi$ -Markov history representation  $\Phi^\pi : H_t^\pi \rightarrow \bar{H}_t^\pi$  with  $\bar{H}_t^\pi = \times_i \text{dom}(\bar{\Theta}_t^{\pi i})$  such that for all  $h_t \in H^\pi$ ,  $\bar{\Theta}_t^\pi \subseteq \bar{\Theta}_t^*$ . Moreover, there exist cases for which  $\bar{\Theta}_t^\pi$  is a proper subset,  $\bar{\Theta}_t^\pi \neq \bar{\Theta}_t^*$ .

*Proof.* The proof follows from Lemma 1. We know that, in general,  $H^\pi \subseteq H$ , and if  $\pi(a'_t | h'_t) = 0$  for at least one pair  $a'_t \in A, h'_t \in H$ , then  $H^\pi \subset H$ . Hence, for every  $\Phi^*$  there is a  $\Phi^\pi$  such that  $\bar{\Theta}_t^\pi \subseteq \bar{\Theta}_t^*$ , and in some cases, when  $H^\pi \subset H$ , we may have  $\bar{\Theta}_t^\pi \neq \bar{\Theta}_t^*$  (e.g., Frozen T-Maze example). □

**Theorem 1.** Let  $\Phi^* : H_t \rightarrow \bar{H}_t^*$  with  $\bar{H}_t^* = \times_i \text{dom}(\bar{\Theta}_t^{*i})$  be a minimal history representation. If, for some  $\pi$ , there is a  $\pi$ -Markov history representation  $\Phi^\pi : H_t^\pi \rightarrow \bar{H}_t^\pi$  with  $\bar{H}_t^\pi = \times_i \text{dom}(\bar{\Theta}_t^{\pi i})$ , such that  $\bar{\Theta}_t^\pi \subset \bar{\Theta}_t^*$  for some  $h_t \in H$ , then  $\Phi^\pi$  is confounded by policy  $\pi$ .

*Proof.* Proof by contradiction. Let us assume that  $\bar{\Theta}_t^\pi \subset \bar{\Theta}_t^*$ , and yet there is no policy confounding. I.e., for all  $h_t, h_{t+1} \in H, a_t \in A$ ,

$$R_h^\pi(\Phi^\pi(h_t), a_t) = R_h^\pi(\text{do}(\Phi^\pi(h_t)), a_t) \quad (1)$$

and

$$\text{Pr}^\pi(\Phi^\pi(h_{t+1}) | \Phi^\pi(h_t), a_t) = \text{Pr}^\pi(\Phi^\pi(h_{t+1}) | \text{do}(\Phi^\pi(h_t)), a_t) \quad (2)$$

First, note that the do-operator implies that the equality must hold for all  $h'_t$  in  $h_t$ 's equivalence class under  $\Phi^\pi$ ,  $h'_t \in \{h_t\}^\Phi = \{h'_t \in H_t : \Phi(h'_t) = \Phi(h_t)\}$ , i.e., not just those  $h'_t$  that are visited under  $\pi$ ,

$$R_h^\pi(\Phi^\pi(h_t), a_t) = R_h^\pi(\text{do}(\Phi^\pi(h_t)), a_t) = \{R(h'_t, a_t)\}_{h'_t \in \{h_t\}^\Phi} \quad (3)$$

which is precisely the first condition in Definition 4,

$$R_h(\Phi^\pi(h_t), a_t) = R_h(h_t, a_t), \quad (4)$$

for all  $h_t \in H$  and  $a_t \in A$ .

Analogously, we have that,

$$\begin{aligned} \Pr^\pi(\Phi^\pi(h_{t+1}) \mid \Phi^\pi(h_t), a_t) &= \Pr^\pi(\Phi^\pi(h_{t+1}) \mid \text{do}(\Phi^\pi(h_t)), a_t) \\ &= \Pr(\Phi^\pi(h_{t+1}) \mid \Phi^\pi(h_t), a_t) \end{aligned} \quad (5)$$

where the second equality reflects that the above must hold independently of  $\pi$ . Hence, we have that for all  $h_t, h_{t+1} \in H$  and  $h'_t \in \{h_t\}^\Phi$ ,

$$\Pr(\Phi^\pi(h_{t+1}) \mid \Phi^\pi(h_t), a_t) = \Pr(\Phi^\pi(h_{t+1}) \mid \Phi^\pi(h'_t), a_t), \quad (6)$$

which means that, for all  $h_t, h_{t+1} \in H$  and  $a_t \in A$ ,

$$\begin{aligned} \Pr(\Phi^\pi(h_{t+1}) \mid \Phi^\pi(h_t), a_t) &= \Pr(\Phi^\pi(h_{t+1}) \mid h_t, a_t) \\ &= \sum_{h'_{t+1} \in \{h_{t+1}\}^{\Phi^\pi}} T_h(h'_{t+1} \mid h_t, a_t), \end{aligned} \quad (7)$$

which is the second condition in Definition 4.

Equations (4) and (7) reveal that if the assumption is true (i.e.,  $\Phi^\pi$  is not confounded by the policy), then  $\Phi^\pi$  is not just  $\pi$ -Markov but actually strictly Markov (Definition 4). However, we know that  $\Phi^*(h_t)$  is the minimal history representation, which contradicts the above statement, since, according to Definition 5, there is no proper subset of  $\bar{\Theta}_t^*$ , for all  $h_t \in H$ , such that the representation remains Markov. Hence,  $\bar{\Theta}_t^\pi \subset \bar{\Theta}_t^*$  implies policy confounding.  $\square$

**Proposition 2.** Let  $\{\bar{\Theta}_t^*\}_{\cup \Phi^*}$  be the union of variables in all possible minimal history representations. There exist cases where, for some  $\pi$ , there is a  $\pi$ -minimal history representation  $\Phi^{\pi*} : H_t^\pi \rightarrow \bar{H}_t^{\pi*}$  with  $\bar{H}_t^{\pi*} = \times_i \text{dom}(\bar{\Theta}_t^{\pi*i})$  such that  $\bar{\Theta}_t^{\pi*} \setminus \{\bar{\Theta}_t^*\}_{\cup \Phi^*} \neq \emptyset$ .

*Proof (sketch).* Consider a deterministic MDP with a deterministic policy. Imagine there exists a variable  $X^1$  that is perfectly correlated with the episode's timestep  $t$ , but that is generally irrelevant to the task. The variable  $X^1$  would constitute in itself a valid  $\pi$ -Markov history representation since it can be used to determine transitions and rewards so long as a deterministic policy is followed. At the same time,  $X^1$  would not enter the minimal Markov history representation because it is useless under stochastic policies. Example 4 below illustrates this situation.  $\square$

## B Example: Watch the Time

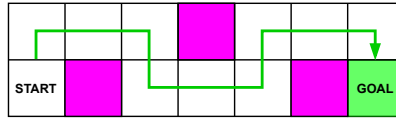


Figure 7: An illustration of the watch-the-time environment.

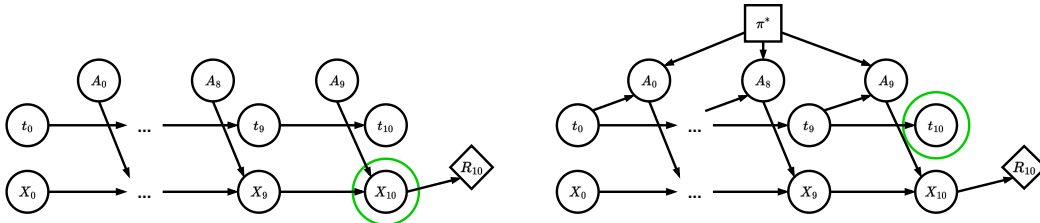


Figure 8: Two DBNs representing the dynamics of the watch-the-time environment, when actions are sampled at random (left), and when they are determined by the optimal policy (right).

**Example 4. (Watch the Time)** This example is inspired by the empirical results of Song et al. (2020). Figure 7 shows a grid world environment. The agent must go from the start cell to the goal cell. The agent must avoid the pink cells; stepping on those yields a  $-0.1$  penalty. There is a  $+1$  reward for reaching the goal. The agent can observe its own location within the maze  $X$  and the current timestep  $t$ . The two diagrams in Figure 8 are DBNs describing the environment dynamics. When actions are considered exogenous random variables (left diagram), the only way to estimate the reward at  $t = 10$  is by looking at the agent’s location. In contrast, when actions are determined by the policy (right diagram),  $t$  becomes a proxy for the agent’s location  $X_{10}$ . This is because the start location and the sequence of actions are fixed. This implies that  $t$  is a perfectly valid  $\pi$ -Markov history representation under  $\pi^*$ . Moreover, as shown by the DBN on the right, the optimal policy may simply rely on  $t$  to determine the optimal action.

## C Further Related Work

**Early evidence of policy confounding** Although to the best of our knowledge, we are the first to bring forward and describe mathematically the idea of policy confounding, a few prior works have reported evidence of particular forms of policy confounding. In their review of the Arcade Learning Environment (ALE; Bellemare et al., 2013), Machado et al. (2018) explain that because the games are fully deterministic (i.e., initial states are fixed and transitions are deterministic), open-loop policies that memorize good action sequences can achieve high scores in ALE. Clearly, this can only occur if the policies themselves are also deterministic. In such cases, policies, acting as confounders, induce a spurious correlation between the past action sequences and the environment states. Similarly, Song et al. (2020) showed, by means of saliency maps, how agents may learn to use irrelevant features of the environment that happen to be correlated with the agent’s progress, such as background clouds or the game timer, as clues for outputting optimal actions. In this case, the policy is again a confounder for all these, a priori irrelevant, features. Zhang et al. (2018b) provide empirical results showing how large neural networks may overfit their training environments and, even when trained on a collection of procedurally generated environments, memorize the optimal action for each observation. Zhang et al. (2018a) shows how, when trained on a small subset of trajectories, agents fail to generalize to a set of test trajectories generated by the same simulator. Lan et al. (2023) report evidence of well-trained agents failing to perform well on Mujoco environments when starting from trajectories (states) that are out of the distribution induced by the agent’s policy. We conceive this as a simple form of policy confounding. Since the Mujoco environments are also deterministic, agents following a fixed policy can memorize the best actions to take for each state instantiation, potentially relying on superfluous features. Hence, they can overfit to unnatural postures that would not occur under different policies. Finally, Nikishin et al. (2022) describe a phenomenon named ‘primacy bias’, which prevents agents trained on poor trajectories from further improving their policies. The authors show that this issue is particularly relevant when training relies heavily on early data coming from a fixed random policy. We hypothesize that one of the causes for this is also policy confounding. The random policy may induce spurious correlations that lead to the formation of rigid history (state) representations that are hard to recover from.

**Generalization** Generalization is a hot topic in machine learning. The promise of a model performing well in contexts other than those encountered during training is undoubtedly appealing. In the realm of reinforcement learning, the majority of research focuses on generalization to environments that, despite sharing a similar structure, differ somewhat from the training environment (Kirk et al., 2023). These differences range from small variations in the transition dynamics (e.g., sim-to-real transfer; Higgins et al., 2017; Tobin et al., 2017; Peng et al., 2018; Zhao et al., 2020), changes in the observations (i.e., modifying irrelevant information, such as noise: Mandlkar et al., 2017; Ornia et al., 2022, or background variables: Zhang et al., 2020; Stone et al., 2021), to alterations in the reward function, resulting in different goals or tasks (Taylor and Stone, 2009; Lazaric, 2012; Muller-Brockhausen et al., 2021). Instead, we focus on the problem of OOT generalization. Keeping the environment unchanged, we aim to ensure that agents perform effectively when confronted with situations that differ from those encountered along their preferred trajectories.

**State abstraction** State abstraction is concerned with removing from the representation all that state information that is irrelevant to the task. In contrast, we are worried about learning representations containing too little information, which can lead to state aliasing. Nonetheless, as argued by

McCallum (1995), state abstraction and state aliasing are two sides of the same coin. That is why we borrowed the mathematical frameworks of state abstraction to describe the phenomenon of policy confounding. Li et al. (2006) provide a taxonomy of the types of state abstraction and how they relate to one another. Givan et al. (2003) introduce the concept of bisimulation, which is equivalent to our definition of Markov history representation (Definition 4) but for states instead of histories. Ferns et al. (2006) proposes a method for measuring the similarity between two states. Castro (2020) notes that this metric is prohibitively expensive and suggests using a relaxed version that computes state similarity relative to a given policy. This is similar to our notion of  $\pi$ -Markov history representation (Definition 7). While the end goal of this metric is to group together states that are similar under a given policy, here we argue that this may lead to poor OOT generalization.

## D Dynamic Bayesian Networks

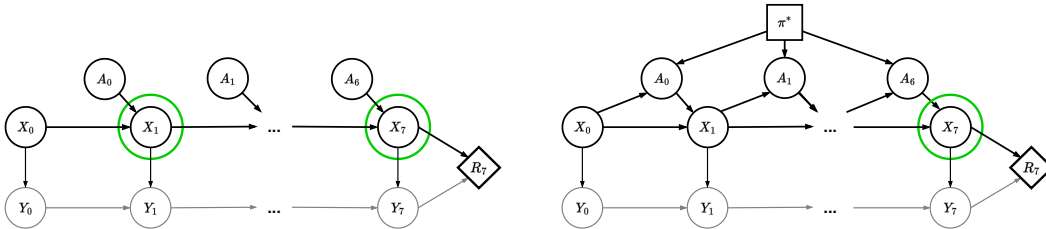


Figure 9: Two DBNs representing the dynamics of the Key2Door environment, when actions are sampled at random (left), and when they are determined by the optimal policy (right). The nodes labeled as  $X$  represent the agent’s location, while the nodes labeled as  $Y$  represent whether or not the key has been collected. The agent can only see  $X$ . Hence, when actions that are sampled are random (left), the agent must remember its past locations to determine the reward  $R_7$ . Note that only  $X_1$  and  $X_7$  are highlighted in the left DBN. However, other variables in  $\langle X_2, \dots, X_6 \rangle$  might be needed, depending on when the key is collected. In contrast, when following the optimal policy, only  $X_7$  is needed. In this second case, knowing the current location is sufficient to determine whether the key has been collected.

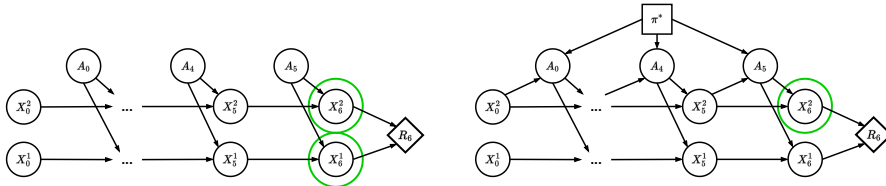


Figure 10: Two DBNs representing the dynamics of the Diversion environment, when actions are sampled at random (left), and when they are determined by the optimal policy (right). The nodes labeled as  $X^1$  indicate the row where the agent is located; the nodes labeled as  $X^2$  indicate the column. We see that when actions are sampled at random, both  $X_6^1$  and  $X_6^2$  are necessary to determine  $R_6$ . However, when actions are determined by the optimal policy,  $X_6^2$  is sufficient, as the agent always stays at the top row.

## E Experimental Results

### E.1 Learned history representations

The results reported in Section 7 show that the OOT generalization problem exists. However, some may still wonder if the underlying reason is truly policy confounding. To confirm this, we compare the outputs of the policy at every state in the Frozen T-Maze when being fed the same histories (observation stack) but two different signals. That is, we permute the variable containing the signal ( $x^1$  in the diagram of Figure 2) and leave the rest of the variables in the observation stack unchanged. We then feed the two versions to the policy network and measure the KL divergence between the two output probabilities. This metric is a proxy for how much the agent attends to the signal in every state.



The heatmaps in Figure 11 show the KL divergences at various points during training (0, 10K, 30K, and 100K timesteps) when the true signal is ‘green’ and we replace it with ‘purple’. We omit the two goal states since no actions are taken there. We see that initially (top left heatmap), the signal has very little influence on the policy (note the scale of the colormap is  $10^{-6}$ ), after 10K steps, the agent learns that the signal is very important when at the top right state (top right heatmap). After this, we start seeing how the influence of the signal at the top right state becomes less strong (bottom left heatmap) until it eventually disappears (bottom right heatmap). In contrast, the influence of the signal at the initial state becomes more and more important, indicating that after taking the first action, the agent ignores the signal and only attends to its own location. The results for the alternative case, purple signal being replaced by green signal, are shown in Figure 12.

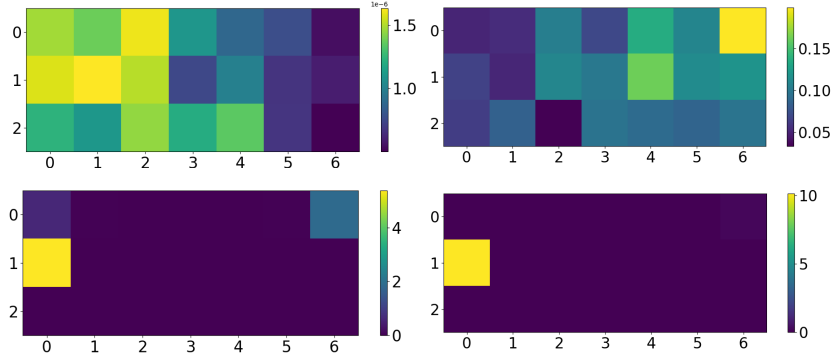


Figure 11: A visualization of the learned history representations. The heatmaps show the KL divergence between the action probabilities when feeding the policy network a stack of the past 10 observations and when feeding the same stack but with the value of the signal being switched from green to purple, after 0 (top left), 10K (top right), 30K (bottom left), and 100K (bottom right) timesteps of training.

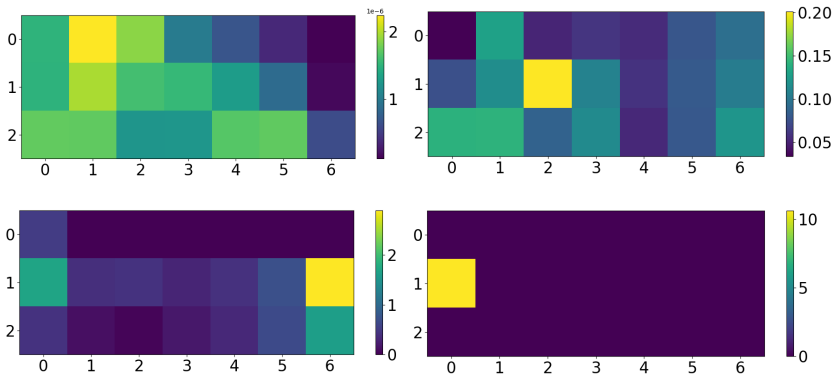


Figure 12: A visualization of the learned history representations. The heatmaps show the KL divergence between the action probabilities when feeding the policy network a stack of the past 10 observations and when feeding the same stack but with the value of the signal being switched from purple to green, after 0 (top left), 10K (top right), 30K (bottom left), and 100K (bottom right) timesteps of training.

## E.2 Buffer size and exploration/domain randomization

Figures 13 and 14 report the results of the experiments described in Section 7 (paragraphs 2 and 3) for Key2Door and Diversion. We see how the buffer size also affects the performance of DQN in the two environments (left plots). We also see that exploration/domain randomization does improve OOT generalization in Diversion but not in Key2Door.

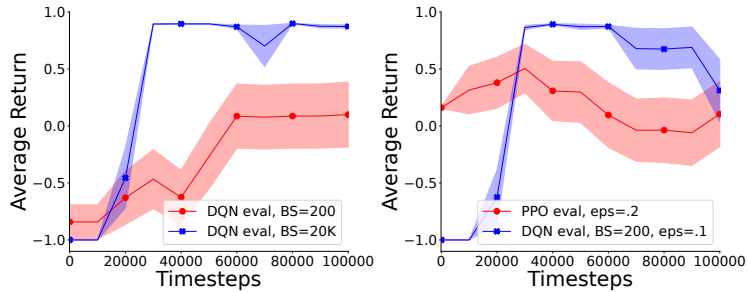


Figure 13: Key2Door. Left: DQN small vs. large buffer sizes. Right: PPO and DQN when adding stochasticity.

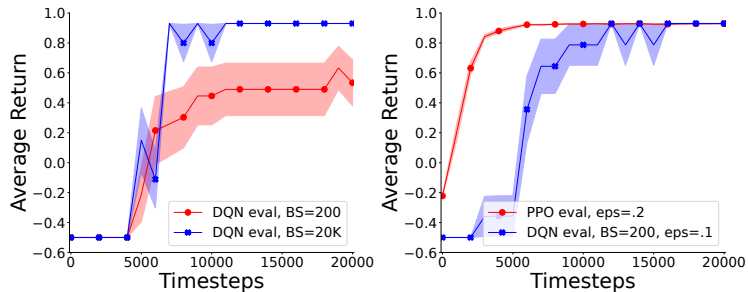


Figure 14: Diversion. Left: DQN small vs. large buffer sizes. Right: PPO and DQN when adding stochasticity.

## F Further Experimental Details

We ran our experiments on an Intel i7-8650U CPU with 8 cores. Agents were trained with Stable Baselines3 (Raffin et al., 2021). Most hyperparameters were set to their default values except for the ones reported in Tables 1 (PPO) and 2 (DQN), which seemed to work better than the default values.

Table 1: PPO hyperparameters.

Rollout steps	128
Batch size	32
Learning rate	2.5e-4
Number epoch	3
Entropy coefficient	1.0e-2
Clip range	0.1
Value coefficient	1
Number Neurons 1st layer	128
Number Neurons 2nd layer	128

Table 2: DQN hyperparameters.

Buffer size	1.0e5
Learning starts	1.0e3
Learning rate	2.5e-4
Batch size	256
Initial exploration bonus	1.0
Final exploration bonus	0.0
Exploration fraction	0.2
Train frequency	5
Number Neurons 1st layer	128
Number Neurons 2nd layer	128