

---

# Object-Aware Regularization for Addressing Causal Confusion in Imitation Learning

---

Jongjin Park<sup>1\*</sup> Younggyo Seo<sup>1\*†</sup> Chang Liu<sup>2</sup> Li Zhao<sup>2</sup>  
Tao Qin<sup>2</sup> Jinwoo Shin<sup>1</sup> Tie-Yan Liu<sup>2</sup>

<sup>1</sup>Korea Advanced Institute of Science and Technology

<sup>2</sup>Microsoft Research Asia

## Abstract

Behavioral cloning has proven to be effective for learning sequential decision-making policies from expert demonstrations. However, behavioral cloning often suffers from the causal confusion problem where a policy relies on the noticeable *effect* of expert actions due to the strong correlation but not the *cause* we desire. This paper presents Object-aware REgularizatiOn (OREO), a simple technique that regularizes an imitation policy in an object-aware manner. Our main idea is to encourage a policy to uniformly attend to all semantic objects, in order to prevent the policy from exploiting nuisance variables strongly correlated with expert actions. To this end, we introduce a two-stage approach: (a) we extract semantic objects from images by utilizing discrete codes from a vector-quantized variational autoencoder, and (b) we randomly drop the units that share the same discrete code together, i.e., masking out semantic objects. Our experiments demonstrate that OREO significantly improves the performance of behavioral cloning, outperforming various other regularization and causality-based methods on a variety of Atari environments and a self-driving CARLA environment. We also show that our method even outperforms inverse reinforcement learning methods trained with a considerable amount of environment interaction.

## 1 Introduction

Imitation learning (IL) holds the promise of learning skills or behaviors directly from expert demonstrations, effectively reducing the need for costly and dangerous environment interaction [21, 45]. Its simplest and effective form is behavioral cloning (BC), which learns a policy by solving a supervised learning problem over state-action pairs from expert demonstrations. While being simple, BC has been successful in a wide range of tasks [4, 7, 31, 33] with careful designs. However, it has been recently evidenced that BC often suffers from the causal confusion problem, where the policy relies on nuisance variables strongly correlated with expert actions, instead of the true *causes* [11, 12, 54].

For example, when we train a BC policy on the Atari Pong environment (see Figure 1a), we observe that a policy relies on nuisance variables in images (i.e., scores) for predicting expert actions, instead of learning the underlying fundamental rule of the environment that experts would have used for making decisions. In particular, Table 1c shows that the policy trained using images with scores struggles to generalize to images with scores masked out (see Figure 1b). However, the policy trained with masked images could generalize to original images with scores, which shows that it successfully learned the rule of the environment. This implies that learning the policy that can identify the true cause of expert actions is important for stable performance at deployment time, where nuisance correlates usually do not hold as in expert demonstrations.

---

\*Equal contribution, in alphabetical order. {jongjin.park, younggyo.seo}@kaist.ac.kr

†This work was done while the author was an intern at Microsoft Research Asia

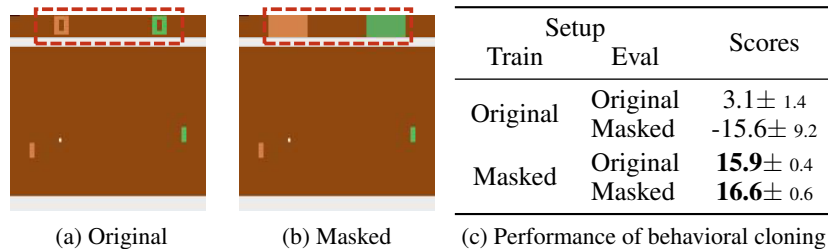


Figure 1: Atari Pong environment with (a) original images and (b) images where scores are masked out. (c) Performance of behavioral cloning (BC) policy trained in Original and Masked environments, averaged over four runs. We observe that the policy trained with original images suffers in both environments, which shows that the policy exploits score information for predicting expert actions, instead of learning the underlying fundamental rule of the environment.

In order to address this causal confusion problem, one can consider causal discovery approaches to deduce the cause-effect relationships from observational data [26, 48]. However, it is difficult to apply these approaches to domains with high-dimensional inputs, as (i) causal discovery from observational data is impossible in general without certain conditions<sup>3</sup> [38], and (ii) these domains usually do not satisfy the assumption that inputs are structured into random variables connected by a causal graph, e.g., objects in images [29, 46]. To address these limitations, de Haan et al. [12] recently proposed a method that learns a policy on top of disentangled representations from a  $\beta$ -VAE encoder [19] with random masking, and infers an optimal causal mask during the environment interaction by querying interactive experts [43] or environment returns. However, given that environment interaction could be dangerous and incur additional costs, we argue that it is important to develop a method for learning the policy robust to causal confusion problem without such a costly environment interaction.

In this paper, we present OREO: **Object-aware REGularizatiOn**, a new regularization technique that addresses the causal confusion problem in imitation learning without environment interaction. The key idea of our method is to regularize a policy to attend uniformly to all semantic objects in images, in order to prevent the policy from exploiting nuisance correlates for predicting expert actions. To this end, we propose to extract semantic objects from raw images by utilizing vector-quantized variational autoencoder (VQ-VAE) [35]. In our experiments, we discover that the units of a feature map corresponding to the objects with similar semantics, e.g., backgrounds, scores, and characters, are mapped into the same or similar discrete codes (see Figure 3). Based upon this observation, we propose to regularize the policy by randomly dropping units that share the same discrete code together throughout training. Namely, our method randomly masks out semantically similar objects, which allows object-aware regularization of the policy.

We highlight the main contributions of this paper below:

- We present OREO, a simple and effective regularization method for addressing the causal confusion problem, and support the effectiveness of OREO with extensive experiments.
- We show that OREO significantly improves the performance of behavioral cloning on confounded Atari environments [5, 12], outperforming various other regularization methods [13, 15, 55, 50] and causality-based methods [12, 47].
- We show that OREO even outperforms inverse reinforcement learning methods trained with a considerable amount of environment interaction [8, 20].

## 2 Related work

**Imitation learning.** Imitation learning (IL) aims to solve complex tasks where learning a policy from scratch is difficult or even impossible, by learning useful skills or behaviors from expert demonstrations [2, 18, 25, 37, 39, 40, 56]. There are two main approaches for IL: inverse reinforcement learning (IRL) methods that find a cost function under which the expert is uniquely optimal [8, 20, 34, 44, 59], and behavioral cloning methods that formulate the IL problem as a supervised learning problem that predicts expert actions from states [3, 4, 7, 31, 33, 40]. Our work employs

<sup>3</sup>de Haan et al. [12] showed that causal discovery methods that depend on faithfulness condition [38] are not applicable to imitation learning setup, as the condition does not hold in environments with nuisance correlates.

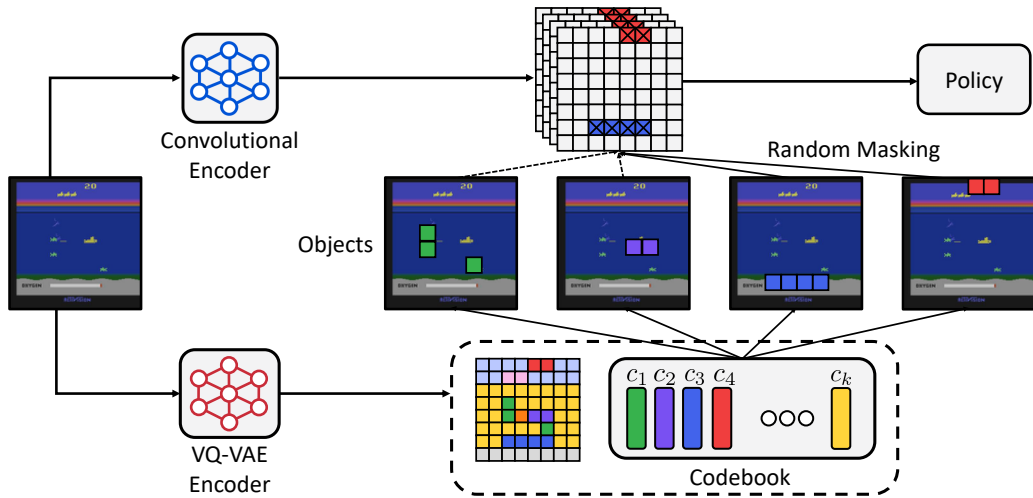


Figure 2: Overview of OREO. We first train a VQ-VAE model that encodes images into discrete codes from a codebook, where each discrete (prototype) representation represents different semantic objects in images. We then regularize a policy by randomly dropping units that share the same discrete code together, i.e., random objects, throughout training.

behavioral cloning as it exhibits the benefit of avoiding costly and dangerous environment interaction, which is crucial for applying imitation learning to real-world scenarios.

**Distributional shift and the causal confusion problem.** Despite its simplicity, BC is known to suffer from the distributional shift, where the state distribution induced by a policy gets different from the training distribution on which the policy was trained. Several approaches have been proposed for learning the policy robust to distributional shift, including interactive IL methods that query experts [42, 43, 51], and regularization techniques [4, 7]. Recently, it has been evidenced that distributional shift leads to the causal confusion problem [4, 12, 54] where a policy exploits the nuisance correlates in states for predicting expert actions. To address this problem, Bansal et al. [4] proposed to randomly drop previous samples from a sequence of samples, and Wen et al. [54] proposed an adversarial training scheme of removing information related to previous actions. The work closest to ours is de Haan et al. [12], which learns a policy with randomly masked disentangled representations and infers the best mask through during environment interaction. Our approach differs in that we regularize the policy to be robust to the causal confusion problem, without any environment interaction.

**Causal discovery from observational data.** Causal discovery aims to discover causal relations among variables by utilizing observational data [38]. Most prior approaches assume that inputs are structured as disentangled variables [6, 16, 26, 36, 48, 47], which often does not hold in domains with high-dimensional inputs, i.e., images. While Lopez-Paz et al. [29] demonstrated the possibility of observational causal discovery from high-dimensional images, combining causal models and representation learning in such domains still remains an open problem [46]. Hence, we instead explore the approach of regularizing a policy that operates on high-dimensional states.

### 3 Method

#### 3.1 Preliminaries

We consider the standard imitation learning (IL) framework where an agent learns to solve a target task from expert demonstrations. Specifically, IL is typically defined in the context of a discrete-time Markov decision process (MDP) [52] without an explicitly-defined reward function, which is defined as a tuple  $(\mathcal{S}, \mathcal{A}, p, \gamma, \rho_0)$ . Here,  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $p(s'|s, a)$  is the transition dynamics,  $\rho_0$  is the initial state distribution, and  $\gamma \in [0, 1)$  is the discount factor. The goal of IL is to learn a policy  $\pi$ , mapping from states to actions, using a set of expert demonstrations  $\mathcal{D} = \{(s_i, a_i)\}_{i=1}^N$ . In our problem setup, an agent cannot interact with the environment, hence it should learn the policy by using only expert demonstrations.

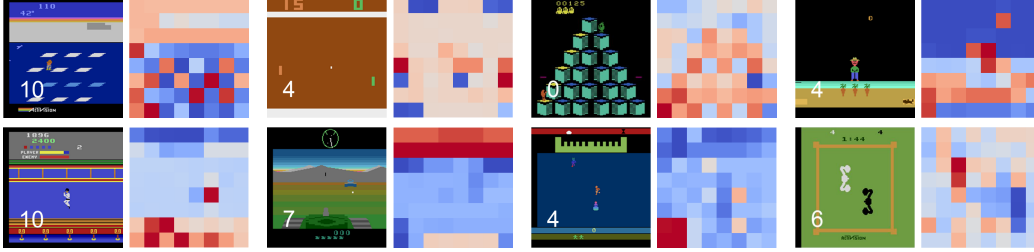


Figure 3: Visualization of the discrete codes from a VQ-VAE model trained on 8 confounded Atari environments, where previous actions are augmented to the images as nuisance variables following the setup in de Haan et al. [12]. The considered environments are Frostbite, Pong, Qbert, Gopher, KungFuMaster, BattleZone, Krull, and Boxing (from left to right, top to bottom). The odd columns show images from environments, and even columns show the corresponding quantized feature maps, respectively. The discrete codes are visualized in 1D using t-SNE [30]. We observe that the units with similar semantics (e.g., the paddles in Pong environment and the carrots in Gopher environment) exhibit similar colors, i.e., mapped into the same or similar discrete codes.

**Behavioral cloning.** Behavioral cloning (BC) reduces an imitation learning problem to the supervised learning problem of training a policy that imitates expert actions. Specifically, we introduce a policy  $\pi$  that maps a state  $s_t$  to an action  $a_t$ , and a convolutional encoder  $f$  that maps a state  $s_t$  to a low-dimensional feature map. Then  $\pi$  and  $f$  are learned by minimizing the negative log-likelihood of expert actions from demonstrations as follows:

$$\mathcal{L}_{\text{BC}}(s_t, a_t) = -\log \pi(a_t | f(s_t)), \quad (1)$$

where  $\pi$  is modeled as a multinomial distribution over actions to handle discrete action spaces.

**Vector quantized variational autoencoder.** The VQ-VAE [35] model consists of an encoder  $g$  that compresses images into discrete latent representations, and a decoder  $d$  that reconstructs images from these discrete representations. Both encoder and decoder share a codebook  $C$  of prototype vectors which are also learned throughout training. Formally, given a state  $s_t$ , the encoder  $g$  encodes  $s_t$  into a feature map  $h_t \in \mathbb{R}^{L \times D}$  that consists of a series of  $L$  latent vectors  $h_{t,i} \in \mathbb{R}^D, i \in \{1, 2, \dots, L\}$ . Then  $h_t = g(s_t)$  is quantized to discrete representations  $e \in \mathbb{R}^{L \times D}$  based on the distance of latent vectors  $h_{t,i}$  to the prototype vectors in the codebook  $C = \{e_k\}_{k=1}^K$  as follows:

$$e_t = (e_{q(t,1)}, e_{q(t,2)}, \dots, e_{q(t,L)}), \quad \text{where } q(t, i) = \underset{k \in [K]}{\operatorname{argmin}} \|h_{t,i} - e_k\|_2, \quad (2)$$

where  $[K]$  is the set  $\{1, \dots, K\}$ . Then the decoder  $d$  learns to reconstruct  $s_t$  from discrete representations  $e_t$ . The VQ-VAE is trained by minimizing the following objective:

$$\mathcal{L}_{\text{VQVAE}}(s_t) = \underbrace{\|s_t - d(e_t)\|_2^2}_{\mathcal{L}_{\text{recon}}} + \underbrace{\|sg[h_t] - e_t\|_2^2}_{\mathcal{L}_{\text{codebook}}} + \underbrace{\beta \cdot \|sg[e_t] - h_t\|_2^2}_{\mathcal{L}_{\text{commit}}}, \quad (3)$$

where the operator  $sg$  refers to a stop-gradient operator,  $\mathcal{L}_{\text{recon}}$  is a reconstruction loss for learning representations useful for reconstructing images,  $\mathcal{L}_{\text{codebook}}$  is a codebook loss to bring codebook representations closer to corresponding encoder outputs  $h$ , and  $\mathcal{L}_{\text{commit}}$  is a commitment loss weighted by  $\beta$  to prevent encoder outputs from fluctuating frequently between different representations.

### 3.2 OREO: Object-aware regularization for behavioral cloning

In this section, we present OREO: **O**bject-aware **R**EGularizati**O**n that regularizes a policy in an object-aware manner to address the causal confusion problem. Our main idea is to encourage the policy to uniformly attend to all semantic objects in images, in order to prevent the policy from exploiting nuisance variables strongly correlated with expert actions. To this end, we introduce a two-stage approach: we first train a VQ-VAE model that encodes images into discrete codes, then learn the policy with our regularization scheme of randomly dropping units that share the same discrete codes (see Figure 2 and Algorithm 1 for the overview and pseudocode of OREO, respectively).

---

**Algorithm 1** Object-aware regularization (OREO)

---

Initialize parameters of encoder  $g$ , decoder  $d$ , codebook  $C$ , policy  $\pi$ .  
**while** not converged **do** // VQ-VAE TRAINING  
  Sample  $\{s_i\}_{i=1}^B \sim \mathcal{D}$ .  
  Update parameters of  $g, d, C$  by minimizing  $\sum_{i=1}^B \mathcal{L}_{\text{VQVAE}}(s_i)$   
**end while**  
Initialize encoder  $f$  with parameters of  $g$ .  
**while** not converged **do** // UPDATE POLICY VIA BEHAVIORAL CLONING  
  Sample  $\{s_i, a_i\}_{i=1}^B \sim \mathcal{D}$   
  Get random masks  $\{M_i\}_{i=1}^B$  in (4)  
  Update parameters of  $f, \pi$  by minimizing  $\sum_{i=1}^B \mathcal{L}_{\text{OREO}}(s_i, a_i, M_i)$  in (5)  
**end while**

---

**Extracting semantic objects.** To regularize a policy in an object-aware manner, we propose to utilize discrete representations from a VQ-VAE model trained by optimizing the objective in (3) with images from expert demonstrations  $\mathcal{D}$ . Our motivation comes from the observation that the units of a feature map corresponding to similar objects are mapped into similar discrete codes (see Figure 3). Then, in order to extract semantic objects from images and utilize them for regularizing the policy, we propose to randomly drop the units of a feature map that share the same discrete code together throughout training. Formally, for each state  $s_t$ , we sample  $K$  binary random variables  $m_k \in \{0, 1\}$ ,  $k = 1, 2, \dots, K$  from a Bernoulli distribution with probability  $1 - p$ , where  $p$  is the drop probability. Then, we construct a mask  $M_t$  by utilizing the discrete representations  $e_t$  in (2) as follows:

$$M_t = (m_{q(t,1)}, m_{q(t,2)}, \dots, m_{q(t,L)}), \quad \text{where } q(t, i) = \underset{k \in [K]}{\operatorname{argmin}} \|h_{t,i} - e_k\|_2. \quad (4)$$

By considering units of a feature map with the same discrete code, we remark that our method can effectively extract semantic objects from high-dimensional images.

**Behavioral cloning with OREO.** Now we propose to utilize our object-aware masking scheme for the regularization of a policy. To this end, we first initialize a convolutional encoder  $f$  with the parameters of a VQ-VAE encoder  $g$ . We empirically find that employing  $f$  as our backbone encoder for  $\pi$  instead of a fixed encoder  $g$  is more effective, as it allows an encoder to learn useful information for predicting actions. Then, we train the policy  $\pi$  by minimizing the following objective:

$$\mathcal{L}_{\text{OREO}}(s_t, a_t, M_t) = -\log \pi(a_t | f(s_t) \odot M_t), \quad (5)$$

where  $\odot$  denotes elementwise product, and the mask  $M_t$  is shared across all channels in a feature map  $f(s_t)$ . Here, our intuition is that our object-aware regularization scheme should be useful for enforcing the policy not to exploit specific objects strongly correlated with expert actions, as the policy should utilize all semantic objects throughout training. Additionally, following Srivastava et al. [50], we scale the masked features by a factor of  $1/(1 - p)$  during the training to ensure the scale of the expected output with masked features to match the scale of outputs at test time.

## 4 Experiments

In this section, we designed our experiments to answer the following questions:

- How does OREO compare to other regularization schemes that randomly drop units from a feature map [15, 50], data augmentation schemes [13, 55], and causality-based methods [12, 47] (see Table 1)?
- How does OREO compare to inverse reinforcement learning methods that learn a policy with environment interaction [8, 20] (see Figure 5)?
- Why is regularization necessary for addressing the causal confusion problem (see Figure 6a), and why is OREO effective for addressing this problem (see Figure 8)?
- Can OREO improve BC using various sizes of expert demonstrations (see Figure 7)?
- Can OREO also address the causal confusion problem when inputs are high-dimensional, complex real-world images (see Table 3)?

Table 1: Performance of policies trained on various confounded Atari environments without environment interaction. OREO achieves the best score on 15 out of 27 environments, and the best median and mean human-normalized score (HNS) over all environments. The results for each environment report the mean of returns averaged over eight runs. We provide standard deviations in Appendix I. CCIL<sup>†</sup> denotes the results without environment interaction.

Environment	BC	Dropout	DropBlock	Cutout	RandomShift	CCIL <sup>†</sup>	CRLR	OREO
Alien	954.1	1003.8	926.4	973.3	806.5	820.0	82.5	<b>1056.2</b>
Amidar	95.8	89.4	110.1	<b>118.7</b>	98.0	74.9	12.0	105.7
Assault	793.8	820.4	815.0	687.6	828.9	683.3	0.0	<b>840.9</b>
Asterix	292.2	313.8	345.4	212.4	135.5	643.2	<b>650.0</b>	180.8
BankHeist	442.1	485.7	508.4	486.1	367.2	<b>653.5</b>	0.0	493.9
BattleZone	11921.2	12457.5	12025.0	11107.5	9180.0	6370.0	1468.8	<b>12700.0</b>
Boxing	18.8	20.3	32.2	20.5	<b>38.3</b>	34.8	-43.0	32.4
Breakout	<b>5.7</b>	5.4	4.8	1.0	2.0	0.5	0.0	4.2
ChopperCommand	874.2	921.4	919.4	1016.1	936.4	760.6	<b>1077.2</b>	977.4
CrazyClimber	45372.9	39501.6	38345.6	44523.2	41924.0	22616.8	112.5	<b>55523.4</b>
DemonAttack	157.2	180.5	167.8	173.1	<b>241.8</b>	171.3	0.0	224.5
Enduro	241.4	250.4	341.8	119.6	316.4	143.1	3.9	<b>522.8</b>
Freeway	32.3	32.4	32.7	32.5	33.0	<b>33.1</b>	21.4	32.7
Frostbite	116.3	124.5	128.2	<b>139.4</b>	121.6	53.3	80.0	129.9
Gopher	1713.9	1819.1	1818.2	1481.0	1995.0	1404.5	0.0	<b>2515.0</b>
Hero	11923.1	14109.7	14711.4	14896.6	12816.0	6567.8	346.2	<b>15219.8</b>
Jamesbond	419.0	451.0	473.8	381.8	428.4	387.2	0.0	<b>502.8</b>
Kangaroo	2781.5	2912.9	3217.1	2824.0	1923.9	1670.5	122.8	<b>3700.2</b>
Krull	3634.3	3892.1	3832.1	3656.4	3788.7	3090.8	0.1	<b>4051.6</b>
KungFuMaster	15074.8	14452.1	15753.0	11405.6	13389.9	13394.9	0.0	<b>18065.6</b>
MsPacman	1432.9	1733.1	1446.4	1711.0	1223.5	1084.2	105.3	<b>1898.4</b>
Pong	3.2	10.2	11.5	6.8	-0.1	-2.7	-21.0	<b>14.2</b>
PrivateEye	2681.8	2599.1	2720.6	2670.6	<b>3969.2</b>	305.3	-1000.0	3124.9
Qbert	5438.4	6469.0	6140.3	5748.6	3921.4	5138.0	125.0	<b>6966.4</b>
RoadRunner	18381.5	21470.9	22265.4	12417.1	16210.0	11834.1	1022.9	<b>24644.2</b>
Seaquest	454.4	471.3	486.8	330.1	<b>1016.8</b>	271.2	172.5	753.1
UpNDown	4221.1	4147.1	<b>4789.2</b>	4159.6	3880.2	2631.1	20.0	4577.9
Median HNS	44.1%	47.4%	49.8%	42.0%	47.6%	36.2%	-1.5%	<b>51.2%</b>
Mean HNS	73.2%	79.0%	91.7%	69.5%	88.1%	71.7%	-45.9%	<b>105.6%</b>

**Environments and datasets.** We evaluate OREO on 27 Atari environments [5], which are selected by following prior works [12, 49]. Following de Haan et al. [12], we consider *confounded* Atari environments, where images are augmented with previous actions (see Figure 4). We utilize a single frame as an input to a policy, to focus on the causal confusion problem from nuisance correlates in the current state<sup>4</sup>. In our experiments, we report two evaluation metrics: average score

from environments and human-normalized score ( $\text{HNS} := \frac{\text{Agent}_{\text{score}} - \text{Random}_{\text{score}}}{\text{Human}_{\text{score}} - \text{Random}_{\text{score}}}$ ), following Mnih et al. [32]. For expert demonstrations, we utilize DQN Replay dataset [1]. As this dataset consists of 50M transitions of each environment collected during the training of a DQN agent [32], we use the last  $N$  trajectories as expert demonstrations. We preprocess input images to grayscale images of  $84 \times 84 \times 1$ , by utilizing Dopamine library [9]. We provide more details in Appendix B.

**Implementation.** We use a single Nvidia P100 GPU and 8 CPU cores for each training run. The training time for OREO is 6 hours on the dataset of size 50000, compared to 3 hours for BC, which is because OREO additionally trains a VQ-VAE model. As for hyperparameter selection, we use the default hyperparameters from previous or similar works [35, 50], i.e., a drop probability of  $p = 0.5$ , a codebook size of  $K = 512$ , and a commitment cost of  $\beta = 0.25$ . We use the same hyperparameters across all environments. We report the results over 8 runs unless specified. Source code and more details on implementation are available in Appendix A and B, respectively.

<sup>4</sup>We refer to Wen et al. [54] for the discussion on the causal confusion problem from stacking states. While we mainly focus on the single frame setup, OREO is also effective on multiple frame setup (see Appendix F).

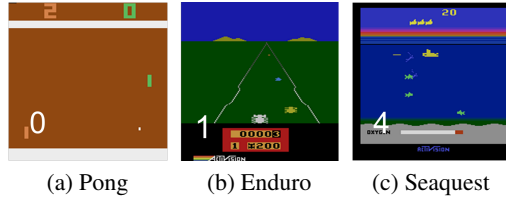


Figure 4: Confounded Atari environments with previous actions (white number in lower left).

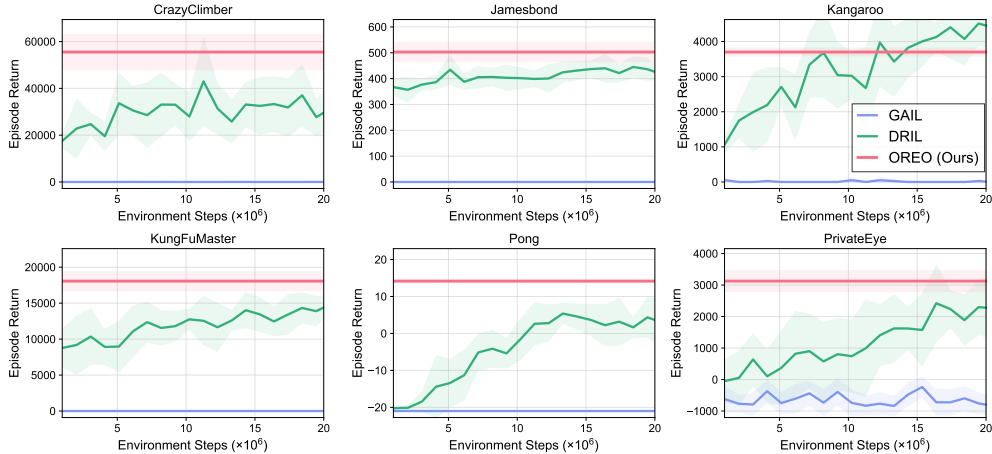


Figure 5: We compare OREO to inverse reinforcement learning methods that require environment interaction for learning a policy, on 6 confounded Atari environments. OREO outperforms baseline methods in most cases, even without using any interaction with environments. The solid line and shaded regions represent the mean and standard deviation, respectively, across eight runs.

**Baselines.** We consider BC as the most basic baseline method. To evaluate the effectiveness of our object-aware regularization scheme, we compare to regularization techniques that drop the randomly sampled units (i.e., Dropout [50]) or randomly sampled blocks (i.e., DropBlock [15]) from the feature map of a convolutional encoder. We also compare to data augmentation schemes, i.e., Cutout [13] that randomly masks out a square patch from images, and RandomShift [55] that randomly shifts pixels of images for regularization. We also consider the method of de Haan et al. [12] that learns a policy on top of disentangled representations from a  $\beta$ -VAE [19] (i.e., CCIL), and an observational causal inference method that estimates the causal contribution of each variable by confounder balancing (i.e., CRLR [47]). We provide the details for baselines in Appendix B and H.

**Comparative evaluation.** Table 1 shows the performance of various methods that learn a policy without environment interaction. OREO significantly improves the performance of BC in most environments, outperforming other regularization techniques. In particular, OREO achieves the mean HNS of 105.6%, while the second-best method, i.e., DropBlock, achieves 91.7%. This demonstrates that our object-aware regularization scheme is indeed effective for addressing the causal confusion problem (see Figure 8 for qualitative experimental results). We found that CCIL without environment interaction does not exhibit strong performance in most environments, possibly due to the difficulty of learning disentangled representations from high-dimensional images [28]. We also provide experimental results for CCIL with environment interaction in Appendix D, where the performance slightly improves but overall trends are similar. We observe that CRLR underperforms in most environments, which shows the difficulty of causal inference from high-dimensional images. We emphasize that OREO also outperforms other baselines in original Atari environments, which implies that our method is also effective for addressing the causal confusion that naturally occurs (see Figure 1). We provide experimental results for the original setup in Appendix C.

**Comparison with inverse reinforcement learning methods.** To demonstrate that OREO can exhibit strong performance without environment interaction, we compare our method to inverse reinforcement learning (IRL) methods that first learn a reward function using expert demonstrations, and train a policy with environment interaction using learned reward function. Specifically, we consider GAIL [20], a method that learns reward function by discriminating expert states from on-policy states during environment interaction; and DRIL [8], one of the strongest IRL methods that utilizes the disagreement between ensemble policies as a reward function. As shown in Figure 5, we observe that OREO exhibits superior performance to GAIL and DRIL on most confounded Atari environments, which are trained with 20M environment steps following the setup in Brantley et al. [8]. While IRL methods might outperform OREO asymptotically with more environment interaction, this result demonstrates that OREO indeed allows for achieving strong performance without interaction. We also found that GAIL exhibits almost zero performance in most environments, which is similar to the observation of previous works that GAIL suffers in environments with high-dimensional image inputs [8, 12, 41]. We remark that OREO can also be applied to IRL methods (see Appendix E for relevant experimental results of DRIL + OREO on confounded Atari environments).

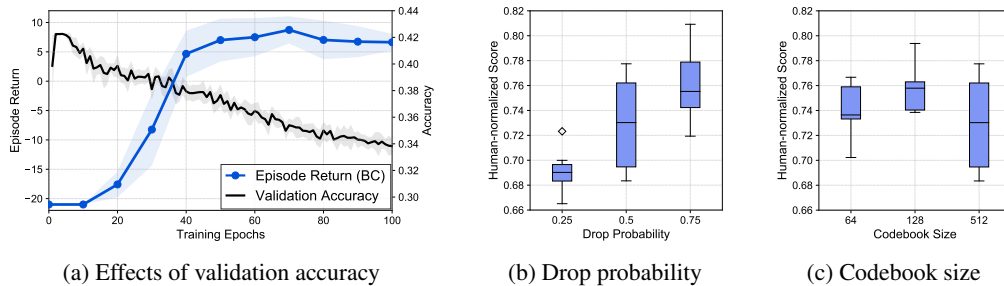


Figure 6: (a) Score and validation accuracy on confounded Pong environment, where validation accuracy is not aligned with the score at test time, which necessitates the use of regularization for addressing the causal confusion problem. We visualize the performance of OREO over 8 confounded Atari environments with varying (b) the drop probability of each code from a codebook and (c) codebook size. Boxplots are drawn using mean human-normalized scores obtained from eight runs.

**Why is regularization necessary in confounded environments?** A simple and widely used approach to address the overfitting problem in supervised learning is a model selection with a validation dataset. To see how this works in our setup, we first introduce a validation dataset consisting of 5 expert demonstrations on confounded Pong environment, and visualize the scores and validation accuracies measured by a policy learned with BC in Figure 6a. We observe that this simple scheme is not helpful for confounded Atari environments, i.e., validation accuracy is not aligned with the score at test time, because the distribution of the validation dataset could be significantly different from the distribution induced by a learned policy. As the evaluation of the policy in environments during training could be dangerous or even impossible, this result implies that regularizing the policy is necessary for successful imitation learning in confounded environments.

**Effects of hyperparameters.** We investigate the effect of two major hyperparameters, i.e.,  $p \in \{0.25, 0.5, 0.75\}$  for the drop probability in (4), and  $K \in \{64, 128, 512\}$  for the codebook size in (2). Figure 6b shows that the performance improves as  $p$  increases, which implies that more strong regularization is effective for addressing the causal confusion problem on confounded Atari environments. Figure 6c shows that too small or large codebook size could be harmful to the performance. We remark that our experiments used default hyperparameters  $p = 0.5$  and  $K = 512$  for reporting the results, so the performance of OREO could be further improved with more tuning.

**Effects of expert demonstration size.** To investigate the effectiveness of OREO with various sizes of expert demonstrations, we evaluate the performance of OREO with a varying number of expert demonstrations  $N \in \{5, 10, 20, 35, 50\}$ . Specifically, we report the mean HNS over 8 confounded Atari environments, which are randomly selected due to the high computation cost of running experiments for all environments. As shown in Figure 7, OREO consistently improves the performance of BC across a wide range of dataset sizes. As for the comparison with other baselines, OREO achieves superior performance to Dropout and DropBlock except for the extreme case of  $N = 5$ , which is because learning a VQ-VAE model with a limited number of data could be unstable. We also observe that DropBlock and OREO consistently outperform Dropout, which supports our intuition that dropping individual units from a feature map is not sufficient for effective regularization to address the causal confusion problem.

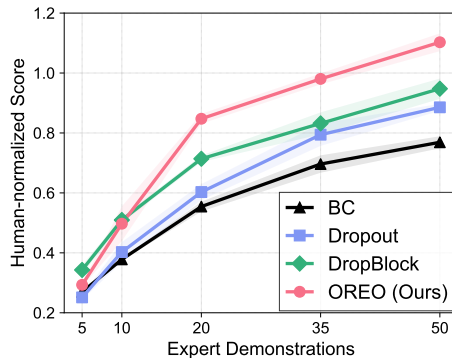


Figure 7: Mean human-normalized score over 8 confounded Atari environments with a varying number of expert demonstrations. The solid line and shaded regions represent the mean and standard deviation, respectively, across four runs.



Table 2: Performance of policies trained on various confounded Atari environments without environment interaction. VQ-VAE + BC learns a BC policy on top of fixed VQ-VAE representations. The results for each environment report the mean and standard deviation of returns over eight runs.

Environment	BC	VQ-VAE + BC	VQ-VAE + Dropout	VQ-VAE + DropBlock	OREO
BankHeist	442.1± 20.7	358.8± 25.8	491.1± 28.9	488.0± 49.7	<b>493.9± 17.6</b>
Enduro	241.4± 28.4	154.6± 10.7	57.1± 12.6	111.2± 16.4	<b>522.8± 29.1</b>
KungFuMaster	15074.8± 275.5	11055.1± 867.2	13323.0± 1390.0	14861.1± 1561.5	<b>18065.6± 1411.5</b>
Pong	3.2± 0.7	3.6± 1.8	10.4± 0.8	13.6± 0.3	<b>14.2± 0.4</b>
PrivateEye	2681.8± 270.2	2255.8± 569.5	390.2± 300.9	746.8± 527.8	<b>3124.9± 349.6</b>
RoadRunner	18381.5± 1519.9	5783.2± 403.6	6633.8± 716.8	7771.1± 843.6	<b>24644.2± 2235.1</b>
Seaquest	454.4± 53.5	344.9± 35.2	325.6± 28.2	396.6± 36.8	<b>753.1± 63.6</b>
UpNDown	4221.1± 214.5	2676.9± 268.9	3310.8± 536.2	4073.9± 760.9	<b>4577.9± 307.6</b>
Median HNS	62.7%	47.9%	45.3%	53.2%	<b>72.9%</b>
Mean HNS	70.8%	41.3%	45.7%	53.0%	<b>100.1%</b>

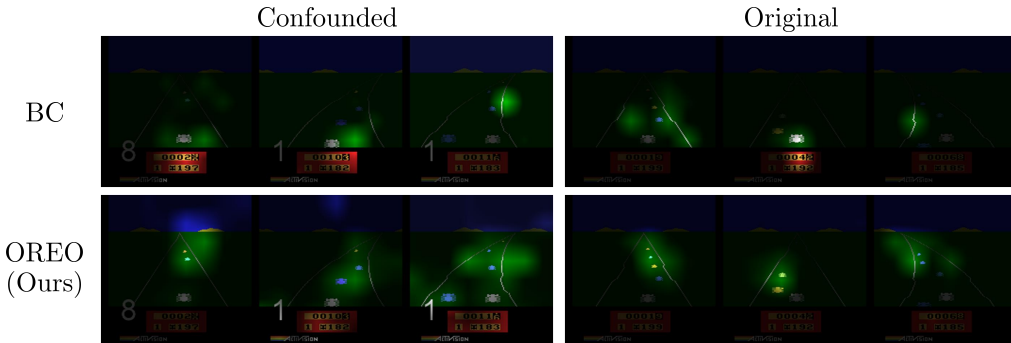


Figure 8: We visualize the spatial attention map from a convolutional encoder trained with BC and OREO in confounded (left) and original (right) Enduro environment. We observe that the encoder trained with OREO attends to important objects from images (e.g., approaching cars), while the encoder trained with BC only attends the small region around a car.

**Contribution of a separate convolutional encoder.** In order to verify the effect of introducing an additional convolutional encoder  $f$  instead of learning a policy on top of the fixed VQ-VAE encoder  $g$  in (5), we provide the experimental results that evaluate VQ-VAE + BC, where a BC policy is learned on top of fixed VQ-VAE representations in Table 2. We first observe that the performance of VQ-VAE + BC performs worse than vanilla BC, which is because fixed VQ-VAE representations learned by reconstruction objective (3) do not contain fine-grained features required for imitating expert actions. Instead, one can see that OREO significantly improves the performance of BC and outperforms all baselines based on fixed VQ-VAE representations, achieving the mean HNS of 100.1% compared to 53.0% of VQ-VAE + DropBlock. This shows that OREO is not the naïve combination of VQ-VAE and BC, but a carefully designed method to exploit the discrete codes from VQ-VAE for object-aware regularization to address the causal confusion problem.

**How does OREO improve the performance of BC?** To understand how OREO improves the performance of BC, we visualize spatial attention maps from the last convolutional layer of a policy encoder in Figure 8. Specifically, following prior works [24, 57], we compute a spatial attention map by averaging the absolute values of a feature map along the channel dimension. We then apply 2-dimensional spatial softmax and multiply the upscaled attention map with images for visualization. We observe that the activations from the encoder trained with OREO are capturing all important objects of the environment (e.g., car that an agent controls and approaching cars), while the activations from BC are missing information of approaching cars by only focusing on the small region around a car and a scoreboard. This shows that our regularization scheme that encourages a policy to uniformly attend to all semantic objects allows for learning the policy that attends to important objects.

**Effectiveness on real-world applications.** To further demonstrate the effectiveness of OREO on real-world applications where inputs are high-dimensional, complex images, we additionally consider a self-driving CARLA environment [14]. Specifically, we train a conditional imitation learning policy

Table 3: Performance of policies trained on 150 expert demonstrations from the CARLA driving dataset, under a weather condition of daytime. The results for each environment report the mean and standard deviation of success rates over four runs. OREO achieves the best success rate on all tasks.

Task	BC	Dropout	DropBlock	OREO
Straight	75.0 $\pm$ 1.7	82.0 $\pm$ 8.3	74.0 $\pm$ 3.5	<b>87.0<math>\pm</math> 4.4</b>
One turn	43.0 $\pm$ 9.1	59.0 $\pm$ 3.3	53.0 $\pm$ 5.2	<b>70.0<math>\pm</math> 7.2</b>
Navigation	16.9 $\pm$ 7.6	30.4 $\pm$ 10.7	21.7 $\pm$ 9.2	<b>35.7<math>\pm</math> 10.2</b>
Navigation w/ dynamic obstacles	18.0 $\pm$ 4.5	26.0 $\pm$ 6.0	19.0 $\pm$ 5.2	<b>30.0<math>\pm</math> 4.5</b>

[10] using 150 expert demonstrations from the dataset [53] consisting of  $200 \times 88 \times 3$  real-world images under a weather condition of daytime. Table 3 shows the average success rate of OREO and baseline methods on four CARLA benchmark tasks, i.e., Straight, One turn, Navigation, and Navigation with dynamics obstacles, where each task consists of 25 different navigation routes. The results show that OREO improves the performance of BC and outperforms other regularization methods, which implies that our object-aware regularization can also be effective on more complex real-world applications.

## 5 Discussion

In this paper, we present OREO, a simple regularization method to address the causal confusion problem in imitation learning. OREO regularizes a policy in an object-aware manner, by randomly dropping the units of a feature map that share the same discrete codes from a VQ-VAE model. Our experimental results demonstrate that OREO improves the performance of behavioral cloning without costly environment interaction, which is crucial for safe and successful imitation learning.

**Limitations.** One limitation of our method is that it is only designed to regularize a policy when inputs are images, and not applicable to state-based environments. However, we still believe that OREO can be a practical solution to the causal confusion problem in various image-based applications, e.g., video games [32], self-driving [7], and robotic manipulation [23]. Another limitation is that we do not deduce the cause-effect relations for addressing the causal confusion problem, but instead regularize the policy to prevent it from exploiting nuisance correlates. However, given that it is an open problem to infer the structured disentangled variables and discover the causal relations among the variables [46], we believe encouraging the policy to attend to *all* semantic objects is a reasonable and promising direction for addressing this problem.

**Potential negative impacts.** Real-world applications of behavioral cloning, e.g., autonomous driving [4], require a large amount of data that often contain sensitive information, therefore raising privacy concerns. As our method is built upon a variational autoencoder, it could be exposed to privacy violation attacks that infer training data information, such as model inversion [58], and membership inference [17]. For example, the facial information of pedestrians may be reconstructed via membership inference attack. To address this vulnerability to privacy violation attacks, a differentially private variational autoencoder would be required for real applications. In addition, pre-training VQ-VAE requires additional computing resources, which might lead to the increased energy cost for learning imitation learning policies. Also, a behavioral cloning policy will imitate whatever demonstrations one specifies. If some bad actions are included in expert demonstrations, the policy would perform dangerous actions to users. For these reasons, in addition to developing algorithms for better performance, it is also important to consider safe adaptation.

## Acknowledgments and Disclosure of Funding

This research was supported by the MSIT (Ministry of Science, ICT), Korea, under the High-Potential Individuals Global Training Program (2020-0-01649) supervised by the IITP (Institute for Information & Communications Technology Planning & Evaluation), Microsoft Research Asia, and Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School

Program(KAIST)). We would like to thank Kimin Lee, Sangwoo Mo, Seonghyeon Park, Sihyun Yu, and anonymous reviewers for providing helpful feedbacks and suggestions in improving our paper.

## References

- [1] Agarwal, Rishabh, Schuurmans, Dale, and Norouzi, Mohammad. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, 2020.
- [2] Aytar, Yusuf, Pfaff, Tobias, Budden, David, Paine, Tom Le, Wang, Ziyu, and de Freitas, Nando. Playing hard exploration games by watching youtube. In *Advances in Neural Information Processing Systems*, 2018.
- [3] Bain, Michael and Sammut, Claude. A framework for behavioural cloning. In *Machine Intelligence 15*, 1996.
- [4] Bansal, Mayank, Krizhevsky, Alex, and Ogale, Abhijit. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. In *Proceedings of Robotics: Science and Systems*, 2019.
- [5] Bellemare, Marc G, Naddaf, Yavar, Veness, Joel, and Bowling, Michael. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- [6] Bengio, Yoshua, Deleu, Tristan, Rahaman, Nasim, Ke, Rosemary, Lachapelle, Sébastien, Bilaniuk, Olexa, Goyal, Anirudh, and Pal, Christopher. A meta-transfer objective for learning to disentangle causal mechanisms. In *International Conference on Learning Representations*, 2020.
- [7] Bojarski, Mariusz, Del Testa, Davide, Dworakowski, Daniel, Firner, Bernhard, Flepp, Beat, Goyal, Prason, Jackel, Lawrence D, Monfort, Mathew, Muller, Urs, Zhang, Jiakai, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [8] Brantley, Kianté, Sun, Wen, and Henaff, Mikael. Disagreement-regularized imitation learning. In *International Conference on Learning Representations*, 2020.
- [9] Castro, Pablo Samuel, Moitra, Subhodeep, Gelada, Carles, Kumar, Saurabh, and Bellemare, Marc G. Dopamine: A research framework for deep reinforcement learning. *arXiv preprint arXiv:1812.06110*, 2018.
- [10] Codevilla, Felipe, Müller, Matthias, López, Antonio, Koltun, Vladlen, and Dosovitskiy, Alexey. End-to-end driving via conditional imitation learning. In *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4693–4700. IEEE, 2018.
- [11] Codevilla, Felipe, Santana, Eder, López, Antonio M, and Gaidon, Adrien. Exploring the limitations of behavior cloning for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.
- [12] de Haan, Pim, Jayaraman, Dinesh, and Levine, Sergey. Causal confusion in imitation learning. In *Advances in Neural Information Processing Systems*, 2019.
- [13] DeVries, Terrance and Taylor, Graham W. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- [14] Dosovitskiy, Alexey, Ros, German, Codevilla, Felipe, Lopez, Antonio, and Koltun, Vladlen. Carla: An open urban driving simulator. In *Conference on robot learning*, pp. 1–16. PMLR, 2017.
- [15] Ghiasi, Golnaz, Lin, Tsung-Yi, and Le, Quoc V. Dropblock: A regularization method for convolutional networks. In *Advances in Neural Information Processing Systems*, 2018.
- [16] Goyal, Anirudh, Lamb, Alex, Hoffmann, Jordan, Sodhani, Shagun, Levine, Sergey, Bengio, Yoshua, and Schölkopf, Bernhard. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*, 2019.

- [17] Hayes, Jamie, Melis, Luca, Danezis, George, and De Cristofaro, Emiliano. Logan: Membership inference attacks against generative models. In *Proceedings on Privacy Enhancing Technologies*, 2017.
- [18] Hester, Todd, Vecerik, Matej, Pietquin, Olivier, Lanctot, Marc, Schaul, Tom, Piot, Bilal, Horgan, Dan, Quan, John, Sendonaris, Andrew, Osband, Ian, et al. Deep q-learning from demonstrations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [19] Higgins, Irina, Matthey, Loic, Pal, Arka, Burgess, Christopher, Glorot, Xavier, Botvinick, Matthew, Mohamed, Shakir, and Lerchner, Alexander. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- [20] Ho, Jonathan and Ermon, Stefano. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, 2016.
- [21] Hussein, Ahmed, Gaber, Mohamed Medhat, Elyan, Eyad, and Jayne, Chrisina. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- [22] Kim, Jaekyeom, Kim, Minjung, Woo, Dongyeon, and Kim, Gunhee. Drop-bottleneck: Learning discrete compressed representation for noise-robust exploration. In *International Conference on Learning Representations*, 2021.
- [23] Kober, Jens and Peters, Jan. Policy search for motor primitives in robotics. *Machine learning*, 84(1-2):171–203, 2011.
- [24] Laskin, Michael, Lee, Kimin, Stooke, Adam, Pinto, Lerrel, Abbeel, Pieter, and Srinivas, Aravind. Reinforcement learning with augmented data. In *Advances in Neural Information Processing Systems*, 2020.
- [25] Le, Hoang M, Yue, Yisong, Carr, Peter, and Lucey, Patrick. Coordinated multi-agent imitation learning. In *International Conference on Machine Learning*, 2017.
- [26] Le, Thuc Duy, Hoang, Tao, Li, Jiuyong, Liu, Lin, Liu, Huawen, and Hu, Shu. A fast pc algorithm for high dimensional causal discovery with multi-core pcs. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 16(5):1483–1495, 2016.
- [27] Liu, Rosanne, Lehman, Joel, Molino, Piero, Such, Felipe Petroski, Frank, Eric, Sergeev, Alex, and Yosinski, Jason. An intriguing failing of convolutional neural networks and the coordconv solution. In *Advances in Neural Information Processing Systems*, 2018.
- [28] Locatello, Francesco, Bauer, Stefan, Lucic, Mario, Raetsch, Gunnar, Gelly, Sylvain, Schölkopf, Bernhard, and Bachem, Olivier. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning*, 2019.
- [29] Lopez-Paz, David, Nishihara, Robert, Chintala, Soumith, Scholkopf, Bernhard, and Bottou, Léon. Discovering causal signals in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [30] Maaten, Laurens van der and Hinton, Geoffrey. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.
- [31] Mahler, Jeffrey and Goldberg, Ken. Learning deep policies for robot bin picking by simulating robust grasping sequences. In *Conference on Robot Learning*, 2017.
- [32] Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David, Rusu, Andrei A, Veness, Joel, Bellemare, Marc G, Graves, Alex, Riedmiller, Martin, Fidjeland, Andreas K, Ostrovski, Georg, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [33] Muller, Urs, Ben, Jan, Cosatto, Eric, Flepp, Beat, and Cun, Yann L. Off-road obstacle avoidance through end-to-end learning. In *Advances in Neural Information Processing Systems*, 2006.
- [34] Ng, Andrew Y, Russell, Stuart J, et al. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning*, 2000.

- [35] Oord, Aaron van den, Vinyals, Oriol, and Kavukcuoglu, Koray. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, 2017.
- [36] Parascandolo, Giambattista, Kilbertus, Niki, Rojas-Carulla, Mateo, and Schölkopf, Bernhard. Learning independent causal mechanisms. In *International Conference on Machine Learning*, 2018.
- [37] Pathak, Deepak, Mahmoudieh, Parsa, Luo, Guanghao, Agrawal, Pulkit, Chen, Dian, Shentu, Yide, Shelhamer, Evan, Malik, Jitendra, Efros, Alexei A, and Darrell, Trevor. Zero-shot visual imitation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018.
- [38] Pearl, Judea. *Causality*. Cambridge University Press, 2009.
- [39] Pohlen, Tobias, Piot, Bilal, Hester, Todd, Azar, Mohammad Gheshlaghi, Horgan, Dan, Budden, David, Barth-Maron, Gabriel, Van Hasselt, Hado, Quan, John, Večerík, Mel, et al. Observe and look further: Achieving consistent performance on atari. *arXiv preprint arXiv:1805.11593*, 2018.
- [40] Pomerleau, Dean. Alvin: An autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems*, 1989.
- [41] Reddy, Siddharth, Dragan, Anca D, and Levine, Sergey. Sqil: Imitation learning via reinforcement learning with sparse rewards. In *International Conference on Learning Representations*, 2020.
- [42] Ross, Stephane and Bagnell, J Andrew. Reinforcement and imitation learning via interactive no-regret learning. In *Advances in Neural Information Processing Systems*, 2014.
- [43] Ross, Stéphane, Gordon, Geoffrey, and Bagnell, Drew. A reduction of imitation learning and structured prediction to no-regret online learning. In *International Conference on Artificial Intelligence and Statistics*, 2011.
- [44] Russell, Stuart. Learning agents for uncertain environments. In *Annual Conference on Computational Learning Theory*, 1998.
- [45] Schaal, Stefan et al. Learning from demonstration. In *Advances in Neural Information Processing Systems*, 1997.
- [46] Schölkopf, Bernhard. Causality for machine learning. *arXiv preprint arXiv:1911.10500*, 2019.
- [47] Shen, Zheyang, Cui, Peng, Kuang, Kun, Li, Bo, and Chen, Peixuan. Causally regularized learning with agnostic data selection bias. In *ACM international conference on Multimedia*, 2018.
- [48] Spirtes, Peter, Glymour, Clark N, Scheines, Richard, and Heckerman, David. *Causation, prediction, and search*. MIT Press, 2000.
- [49] Srinivas, Aravind, Laskin, Michael, and Abbeel, Pieter. Curl: Contrastive unsupervised representations for reinforcement learning. In *International Conference on Machine Learning*, 2020.
- [50] Srivastava, Nitish, Hinton, Geoffrey, Krizhevsky, Alex, Sutskever, Ilya, and Salakhutdinov, Ruslan. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [51] Sun, Wen, Venkatraman, Arun, Gordon, Geoffrey J, Boots, Byron, and Bagnell, J Andrew. Deeply aggravated: Differentiable imitation learning for sequential prediction. In *International Conference on Machine Learning*, 2017.
- [52] Sutton, Richard S and Barto, Andrew G. *Reinforcement learning: An introduction*. MIT Press, 2018.

- [53] Tai, Lei, Yun, Peng, Chen, Yuying, Liu, Congcong, Ye, Haoyang, and Liu, Ming. Visual-based autonomous driving deployment from a stochastic and uncertainty-aware perspective. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2622–2628. IEEE, 2019.
- [54] Wen, Chuan, Lin, Jierui, Darrell, Trevor, Jayaraman, Dinesh, and Gao, Yang. Fighting copycat agents in behavioral cloning from observation histories. In *Advances in Neural Information Processing Systems*, 2020.
- [55] Yarats, Denis, Kostrikov, Ilya, and Fergus, Rob. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International Conference on Learning Representations*, 2021.
- [56] Ye, Gu and Alterovitz, Ron. Guided motion planning. In *Robotics Research*, pp. 291–307. Springer, 2017.
- [57] Zagoruyko, Sergey and Komodakis, Nikos. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017.
- [58] Zhang, Yuheng, Jia, Ruoxi, Pei, Hengzhi, Wang, Wenxiao, Li, Bo, and Song, Dawn. The secret revealer: Generative model-inversion attacks against deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [59] Ziebart, Brian D, Maas, Andrew L, Bagnell, J Andrew, and Dey, Anind K. Maximum entropy inverse reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2008.

# Appendix

## A Source codes

Source codes for reproducing our experimental results are available at <https://github.com/alinlab/oreo>.

## B Details on Atari experiments

### B.1 Experimental setup

**Environments and datasets.** We utilize DQN Replay dataset<sup>5</sup> [1] for expert demonstrations on 27 Atari environments [5]. To encourage the size of the dataset to be consistent across multiple environments, we use the number of expert demonstrations  $N \in \{20, 50\}$ . We provide the size of a dataset for each environment in Table 4. We process input images to grayscale images of  $84 \times 84 \times 1$ , by utilizing Dopamine library<sup>6</sup> [9]. Following de Haan et al. [12], we consider *confounded* Atari environments, where images are augmented with previous actions (see Figure 4). We provide source codes for loading images from the dataset, preprocessing images, and augmenting numbers to the images in Section A. For experiments with selected environments in Figure 7, we randomly chose 8 confounded Atari environments, i.e., BankHeist, Enduro, KungFuMaster, Pong, PrivateEye, RoadRunner, Seaquest, and UpNDown, due to the high computational cost of considering all environments.

**Evaluation.** (a) For all experimental results without environment interaction, we train a policy for 1000 epochs without early stopping based on validation accuracy (see Figure 6a for how early stopping is not effective in our setup), and report the final performance of the trained policy. Specifically, we average the scores over 100 episodes evaluated on confounded environments for each random seed. (b) For all experimental results with inverse reinforcement learning methods that require environment interaction (i.e., GAIL [20] and DRIL [8]), we evaluate a policy over 10 episodes every 1M environment step during the training.

Table 4: Dataset size of each Atari environment.

Environment	$N$	Data	Environment	$N$	Data	Environment	$N$	Data
Alien	50	53165	CrazyClimber	20	83557	Krull	50	65701
Amidar	20	57155	DemonAttack	20	47727	KungFuMaster	20	57235
Assault	50	58868	Enduro	20	169767	MsPacman	50	64305
Asterix	20	68126	Freeway	20	41020	Pong	20	41402
BankHeist	50	58516	Frostbite	50	24043	PrivateEye	20	54020
BattleZone	50	83061	Gopher	20	44011	Qbert	50	59379
Boxing	50	47170	Hero	50	68903	RoadRunner	50	60546
Breakout	50	63799	Jamesbond	20	33659	Seaquest	20	37682
ChopperCommand	50	35262	Kangaroo	20	45898	UpNDown	50	74348

### B.2 Implementation details

#### Implementation details for OREO.

- **VQ-VAE training.** We use the publicly available implementation of VQ-VAE<sup>7</sup> modified to make it work with images of size  $84 \times 84 \times 1$ . Specifically, The encoder consists of four convolutional layers, three with stride 2 and kernel size  $4 \times 4$  and one with stride 1 and kernel size  $3 \times 3$ , followed by 2 residual  $3 \times 3$  blocks (implemented as ReLU,  $3 \times 3$  Conv, ReLU,  $1 \times 1$  conv), all having 256 hidden units. The decoder similarly has 2 residual  $3 \times 3$  blocks, followed by four transposed convolution layers, one with stride 1 and kernel size  $3 \times 3$ , and three with stride 2

<sup>5</sup><https://research.google/tools/datasets/dqn-replay>

<sup>6</sup><https://github.com/google/dopamine>

<sup>7</sup><https://github.com/zalandoresearch/pytorch-vq-vae>

and kernel size  $4 \times 4$ . For training, we train a VQ-VAE model for 1000 epochs with a batch size of 1024. We use Adam optimizer with the learning rate of  $3e-4$ . As for the hyperparameters of VQ-VAE, we use a codebook size of  $K = 512$ , and a commitment cost of  $\beta = 0.25$ , following the original implementation.

- **Behavioral cloning with OREO.** For efficient implementation of OREO, we first compute quantized discrete codes of all images in datasets with pre-trained VQ-VAE, instead of processing every image through VQ-VAE encoder during the training. Then we utilize stored discrete codes for obtaining random masks for training a policy. We find that generating multiple random masks for each image and aggregating the loss computed with each mask marginally improves the performance, by providing more diverse features to the policy. In our experiments, we generate 5 random masks during training. We train a policy for 1000 epochs with the batch size of 1024, and use Adam optimizer with the learning rate of  $3e-4$ .

### Implementation details for regularization and causality-based methods.

- **Behavioral cloning.** We train a BC policy by optimizing the objective in (1) using states and actions from expert demonstrations. Note that other regularization baselines are based on BC.
- **Dropout.** Dropout [50] is a regularization technique that drops units of a feature map from a convolutional encoder. Specifically, for all units of a feature map, Dropout samples binary random variables from a Bernoulli distribution with probability  $1 - p$ , and apply the randomly sampled masks throughout training. We use `nn.Dropout` from PyTorch<sup>8</sup> library with  $p = 0.5$ .
- **DropBlock.** DropBlock [15] is a regularization technique that drops units in a contiguous region of a feature map, i.e., blocks, with the default hyperparameters of  $p = 0.3$  and the block size of 3. We use the publicly available implementation of DropBlock<sup>9</sup> for our experiments. Following this original implementation, we linearly increase  $p$  from 0 to the target value during training.
- **Cutout.** Cutout [13] randomly masks out a square patch from images. We randomly sampled the size of the patch from  $10 \times 10$  to  $30 \times 30$ , by using `RandomErasing` from Kornia<sup>10</sup> library.
- **RandomShift.** RandomShift [55] is a regularization technique that *shifts* images by randomly sampled pixels. Specifically, it pads each side of an image by 4 pixels with boundary pixels and performs random crop of size  $84 \times 84$ . We implemented RandomShift by following the publicly available implementation<sup>11</sup> from the authors.
- **CCIL.** CCIL (named after Causal Confusion in Imitation Learning; [12]) is an interventional causal discovery method that first (i) learns disentangled representations from  $\beta$ -VAE [19] and (ii) infers the causal graph during environment interaction. As publicly available implementation<sup>12</sup> only contains source code that works on the low-dimensional MountainCar environment, we faithfully reproduce the method and report the results. Specifically, we employ `CoordConv`<sup>13</sup> [27] for both the encoder and decoder architectures of  $\beta$ -VAE. We find that prediction accuracy of a policy trained using a fixed  $\beta$ -VAE does not improve over chance level accuracy, possibly because a reconstruction task is not sufficient for learning representations that capture the information required for predicting actions. Hence, we additionally introduce an action prediction task when training a  $\beta$ -VAE, which we find crucial for improving the accuracy over chance level accuracy.
- **CRLR.** As CRLR requires inputs to be binary values, we develop and compare to the categorical version of CRLR that works on top of VQ-VAE discrete codes (see Appendix H).

**Implementation details for inverse reinforcement learning methods.** For all inverse reinforcement learning (IRL) methods, we use the publicly available implementation (<https://github.com/xkianteb/dril>) for reporting the results, with additional modification to original source code to train and evaluate a policy on confounded Atari environments.

<sup>8</sup><https://pytorch.org>

<sup>9</sup><https://github.com/miguelvr/dropblock>

<sup>10</sup><https://github.com/kornia/kornia>

<sup>11</sup><https://github.com/denisyarats/drq>

<sup>12</sup><https://github.com/pimdh/causal-confusion>

<sup>13</sup><https://github.com/walsvid/CoordConv>



- **GAIL.** GAIL [20] is an IRL method that learns a discriminator network that distinguishes expert states from states visited by the current policy, and utilizes the negative output of the discriminator as a reward signal for learning RL agents during environment interaction.
- **DRIL.** DRIL [8] is an IRL method that learns an ensemble of behavioral cloning policies and utilizes the disagreement (i.e., variance) between the predictions of ensemble policies as a cost signal (the negative of reward signal) for learning RL agents during environment interaction.

## C Comparative evaluation on original Atari environments

Table 5 shows the performance of various methods which do not use environment interaction, on original Atari environments. We observe that OREO significantly improves behavioral cloning, also outperforming baseline methods. In particular, OREO achieves the mean HNS of 114.9%, while the second-best method, i.e., DropBlock, achieves 99.0%. This demonstrates that our object-aware regularization scheme is also effective for addressing the causal confusion that naturally occurs in the dataset (see Figure 1).

Table 5: Performance of policies trained on various original Atari environments without environment interaction. OREO achieves the best score on 14 out of 27 environments, and the best median and mean human-normalized score (HNS) over all environments. The results for each environment report the mean of returns averaged over eight runs. CCIL<sup>†</sup> denotes the results without environment interaction.

Environment	BC	Dropout	DropBlock	Cutout	RandomShift	CCIL <sup>†</sup>	CRLR	OREO
Alien	986.5	1117.2	1094.8	1104.4	863.5	1050.4	100.0	<b>1222.2</b>
Amidar	90.8	81.6	113.5	125.0	78.2	78.6	12.0	<b>130.5</b>
Assault	816.8	901.1	829.9	694.1	848.7	755.5	0.0	<b>905.2</b>
Asterix	249.0	176.6	252.2	195.0	99.1	314.1	<b>592.5</b>	212.5
BankHeist	399.0	476.6	471.2	442.5	354.8	<b>606.1</b>	0.0	448.4
BattleZone	10933.8	11621.2	<b>12067.5</b>	10641.2	8748.8	11191.2	5615.0	11703.8
Boxing	21.8	25.7	32.1	21.2	35.8	34.2	-43.0	<b>39.9</b>
Breakout	<b>6.4</b>	2.9	6.0	3.1	4.4	2.1	0.0	5.4
ChopperCommand	1163.0	1162.0	1161.8	1183.9	1026.2	1027.2	1070.2	<b>1282.9</b>
CrazyClimber	54142.2	54965.4	55854.0	47456.4	60465.9	39015.2	885.5	<b>69380.1</b>
DemonAttack	238.8	<b>359.3</b>	225.6	217.8	294.8	194.6	22.7	0.0
Enduro	226.2	304.6	359.1	132.9	282.2	182.8	0.8	<b>514.4</b>
Freeway	32.3	32.6	32.6	32.8	33.0	<b>33.1</b>	21.4	32.9
Frostbite	153.6	149.2	<b>165.7</b>	135.2	133.2	96.7	78.1	152.7
Gopher	1874.4	2220.4	2040.5	1588.2	1456.2	1301.9	0.0	<b>2903.9</b>
Hero	15100.4	15994.4	17058.6	15971.8	14867.2	<b>17487.6</b>	0.0	16370.3
Jamesbond	447.6	492.3	481.9	418.9	452.1	460.4	0.0	<b>527.9</b>
Kangaroo	3162.8	2860.4	<b>3638.6</b>	3242.6	2202.1	2938.1	0.0	3602.9
Krull	4447.9	<b>4764.7</b>	4526.5	4270.6	4611.6	4247.1	0.0	4633.6
KungFuMaster	12900.6	14994.5	14819.0	9956.9	11698.0	12876.9	0.0	<b>16955.5</b>
MsPacman	1921.9	2022.6	2151.7	1949.7	1046.3	1160.6	70.0	<b>2263.8</b>
Pong	3.7	10.0	11.6	7.8	0.8	-19.8	-21.0	<b>12.5</b>
PrivateEye	3035.4	3396.3	3057.6	3092.2	<b>3578.9</b>	1016.4	-1000.0	3162.6
Qbert	5925.4	<b>6363.1</b>	5904.3	6174.8	4100.1	5056.3	125.0	5763.4
RoadRunner	18010.1	20137.8	22522.5	12698.9	15615.4	18985.2	1528.6	<b>27303.9</b>
Seaquest	527.5	644.4	622.3	376.6	<b>948.0</b>	402.4	169.8	921.0
UpNDown	3782.1	3504.3	3886.4	3675.9	3500.4	3062.3	20.0	<b>4186.8</b>
Median HNS	46.7%	53.3%	47.7%	42.9%	47.3%	36.8%	-1.5%	<b>53.6%</b>
Mean HNS	82.0%	91.5%	99.0%	75.0%	91.7%	85.4%	-45.4%	<b>114.9%</b>

## D CCIL with environment interaction

In this section, we compare CCIL with environment interaction, which employs targeted intervention during environment interaction. Specifically, CCIL infers a causal mask over disentangled latent variables from  $\beta$ -VAE, by utilizing the returns from environments. As shown in Figure 9, the performance of CCIL improves during environment interaction of 100 episodes, but OREO still exhibits superior performance to CCIL on most confounded Atari environments. This again demonstrates the difficulty of learning disentangled representations from high-dimensional images [28].

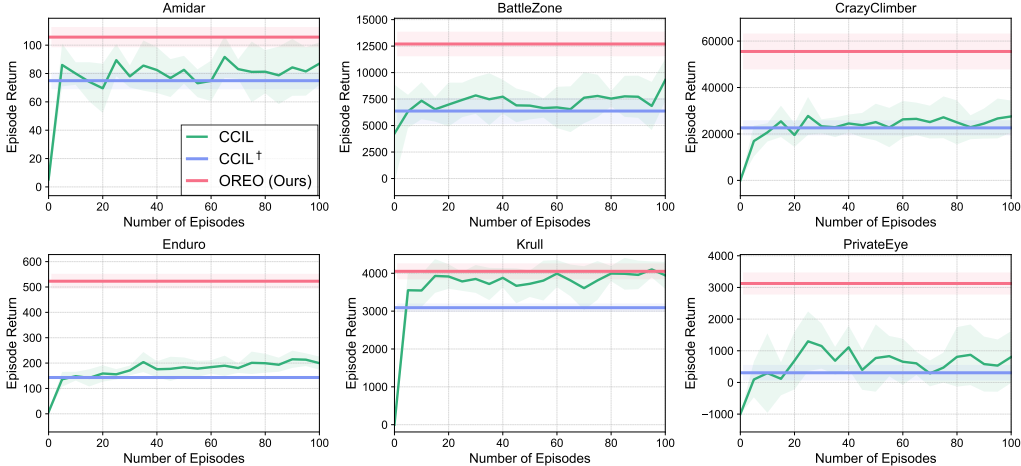


Figure 9: We compare OREO to CCIL with environment interaction, on 6 confounded Atari environments.  $CCIL^\dagger$  denotes the results without environment interaction. The solid line and shaded regions represent the mean and standard deviation, respectively, across eight runs. OREO still outperforms CCIL in most cases, although environment interaction slightly improves the performance of  $CCIL^\dagger$ .

## E Applying OREO to inverse reinforcement learning

We investigate the possibility of applying OREO to other IL methods. While there could be various approaches to utilize the proposed approach for utilizing our regularization scheme for IL, we consider a straightforward application of OREO to a state-of-the-art IL method, i.e., DRIL [8]. Specifically, we apply OREO to the components of DRIL which involves behavioral cloning, i.e., initializing a BC policy and computing rewards with an ensemble of BC policies. In Figure 10, we observe that DRIL + OREO improves the sample-efficiency of DRIL since OREO enables us to learn high-quality BC policies that also result in high-quality reward signals which boosts sample-efficiency. We remark that these results show that IRL methods can also suffer from the causal confusion problem, and a proper regularization scheme can improve the performance by addressing the confusion problem.

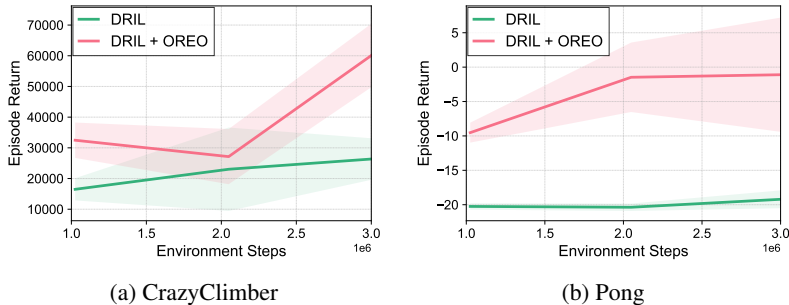


Figure 10: We apply OREO to the inverse reinforcement learning method (i.e., DRIL [8]) and observe that OREO improves the sample-efficiency of DRIL on confounded CrazyClimber and Pong environments. The solid line and shaded regions represent the mean and standard deviation, respectively, across four runs.

## F OREO with a sequence of observations

A natural extension of OREO is to apply our regularization scheme to address the causal confusion problem from a sequence of observations [4, 54]. By extracting semantic objects with the same discrete code from consecutive images and dropping the codes from all images, OREO can regularize the policy consistently over multiple images. In this section, we investigate the effectiveness of OREO on such setup by providing additional experimental results on confounded environments where inputs are four stacked observations. Specifically, we mask the features that correspond to the same discrete codes from each observation, and utilize the aggregated masked features for policy learning. In Table 6, we observe that OREO significantly improves the performance of BC, which shows that OREO can also be effective on this setup by regularizing the policy consistently over multiple frames.

Table 6: Performance of policies trained with four stacked observations on 8 confounded Atari environments. The results for each environment report the mean and standard deviation of returns over four runs.

Environment	BC	Dropout	DropBlock	OREO
BankHeist	448.6± 17.8	477.6± 36.4	466.2± 17.5	<b>538.8± 13.9</b>
Enduro	167.8± 31.7	253.1± 21.1	172.6± 18.4	<b>426.0± 18.1</b>
KungFuMaster	13523.5± 831.7	15041.0± 1011.8	14859.2± 1242.6	<b>18375.2± 1055.3</b>
Pong	4.8± 0.9	8.2± 0.2	9.5± 0.4	<b>12.2± 0.4</b>
PrivateEye	2349.4± 253.1	2173.8± 168.3	<b>2611.4± 476.8</b>	2580.7± 484.2
RoadRunner	15189.5± 1829.0	16574.0± 2799.3	16901.0± 1790.1	<b>18726.2± 876.5</b>
Seaquest	353.4± 11.8	351.4± 28.6	315.3± 17.7	<b>393.2± 19.7</b>
UpNDown	4075.5± 165.6	4306.6± 216.9	4448.9± 450.5	<b>5193.7± 513.5</b>
Median HNS	56.6%	65.1%	59.3%	<b>76.7%</b>
Mean HNS	62.3%	71.0%	68.8%	<b>87.4%</b>

## G Comparison with DropBottleneck

In this section, we compare OREO with DropBottleneck (DB; [22]), which is a dropout-based method that drops features from input variable  $X$  redundant for predicting target variable  $Y$ . While this method was successfully applied to remove the dynamics-irrelevant information such as noises by setting input variable  $X$  and target variable  $Y$  to two consecutive states, we remark that removing task-irrelevant information cannot be an effective recipe for addressing the causal confusion problem. This is because the causal confusion comes from the difficulty of identifying the true cause of expert actions when both confounders and the causes are strongly correlated with expert actions, i.e., they are both task-relevant information. To support this, we provide experimental results where we jointly optimize DB objective when training a BC policy, i.e., setting the target variable  $Y$  to expert actions (denoted as **DB (Y=action)**) in Table 7. In addition, following the original setup in [22], we also provide experimental results where input and target variables are consecutive two states (denoted as **DB (Y=state)**) in Table 8. We observe that **DB (Y=action)** shows comparable performance to OREO in some environments (e.g., CrazyClimber), but OREO still significantly outperforms the suggested baseline in most environments (e.g., Alien, KungFuMaster, and Pong). **DB (Y=state)** performs no better than BC in most environments except for CrazyClimber. These results show that removing dynamics-irrelevant information might not be enough for addressing the causal confusion problem.

Table 7: The results for each environment report the mean and standard deviation of returns over four (DB with expert action) or eight (others) runs. As for the scale of compression term  $\beta$  in DB, we choose a better hyperparameter from an array of [0.001, 0.0001].

Environments	BC	Dropout	DropBlock	DB (Y=action)	OREO
Alien	954.1 ± 83.9	1003.8 ± 53.6	926.4 ± 70.5	994.5 ± 85.6	<b>1056.2 ± 61.6</b>
CrazyClimber	45372.9 ± 5508.9	39501.6 ± 6499.3	38345.6 ± 7190.8	<b>60996.8 ± 7943.5</b>	55523.4 ± 7722.2
KungFuMaster	15074.8 ± 275.5	14452.1 ± 865.4	15753.0 ± 1265.2	15139.5 ± 867.4	<b>18065.6 ± 1411.5</b>
Pong	3.2 ± 0.7	10.2 ± 1.3	11.5 ± 1.3	8.2 ± 0.4	<b>14.2 ± 0.4</b>

Table 8: The results for each environment report the mean and standard deviation of returns over four (DB with consecutive state) or eight (others) runs. As for the scale of compression term  $\beta$  in DB, we choose a better hyperparameter from an array of [0.001, 0.0001].

Environments	BC	Dropout	DropBlock	DB (Y=state)	OREO
Alien	954.1 $\pm$ 83.9	1003.8 $\pm$ 53.6	926.4 $\pm$ 70.5	896.4 $\pm$ 10.7	<b>1056.2 <math>\pm</math> 61.6</b>
CrazyClimber	45372.9 $\pm$ 5508.9	39501.6 $\pm$ 6499.3	38345.6 $\pm$ 7190.8	<b>60111.5 <math>\pm</math> 5597.8</b>	55523.4 $\pm$ 7722.2
KungFuMaster	15074.8 $\pm$ 275.5	14452.1 $\pm$ 865.4	15753.0 $\pm$ 1265.2	15014.3 $\pm$ 1056.2	<b>18065.6 <math>\pm</math> 1411.5</b>
Pong	3.2 $\pm$ 0.7	10.2 $\pm$ 1.3	11.5 $\pm$ 1.3	3.5 $\pm$ 2.1	<b>14.2 <math>\pm</math> 0.4</b>

## H A categorical version of CRLR

In this section, we provide a categorical version of Causally Regularized Logistic Regression (CRLR [47]) method. We first formulate the problem setup and briefly introduce some background on CRLR. Given the training data  $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^N$ , where  $x \in \mathbb{R}^d$  represents the features and  $y$  represents labels, the causal classification task targets to jointly identify the causal contribution  $\beta \in \mathbb{R}^d$  for all features and learn a classifier  $f(\cdot)$  based on  $\beta$ . As we have no prior knowledge of the causal structure, a reasonable way to adapt causal inference into the classification task is to regard each feature  $x_j$  as a treated variable, and all the remaining features  $x_{-j} = x \setminus x_j$  as confounding variables, i.e., confounders. To safely estimate the causal contribution of a given feature  $x_j$ , one has to remove the confounding bias induced by the different distributions of confounders  $x_{-j}$  between the treated and control groups. CRLR finds optimal sample weights to balance the distribution of the treated and control group for any treated variable, under an assumption of binary features. To this end, CRLR learns those sample weights by minimizing a causal regularizer as follows:

$$\min_{\{w^{(i)}\}} \sum_j \left\| \frac{\sum_{i_0 \in I_{j,0}} w^{(i_0)} x_{-j}^{(i_0)}}{\sum_{i_0 \in I_{j,0}} w^{(i_0)}} - \frac{\sum_{i_1 \in I_{j,1}} w^{(i_1)} x_{-j}^{(i_1)}}{\sum_{i_1 \in I_{j,1}} w^{(i_1)}} \right\|_2^2,$$

where  $w^{(i)}$  is the sample weight for  $x^{(i)}$ , and  $I_{j,c}$  denotes  $\{i \mid x_j^{(i)} = c\}$ . The original version of CRLR is built upon the binary features, however, it can be naturally extended to a categorical version, by computing the confounder balancing term for any pair of categorical variables. We convert given categorical features  $x = [x_1, \dots, x_d]$  to one-hot encoded binary features  $\mathbf{x}$ , i.e.,  $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_d]$  where  $\mathbf{x}_i$  is an one-hot encoded version of each feature  $x_i$ . We denote  $\mathbf{x}_{-j} = [\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, 0, \mathbf{x}_{j+1}, \dots, \mathbf{x}_d]$  as confounding variables of these one-hot features. Then, a categorical version of the causal regularizer is computed as follows:

$$\min_{\{w^{(i)}\}} \sum_j \sum_{\substack{c_1 < c_2 \\ \in \{c \in [K] \mid |I_{j,c}| > 0\}}} \left\| \frac{\sum_{i_1 \in I_{j,c_1}} w^{(i_1)} \mathbf{x}_{-j}^{(i_1)}}{\sum_{i_1 \in I_{j,c_1}} w^{(i_1)}} - \frac{\sum_{i_2 \in I_{j,c_2}} w^{(i_2)} \mathbf{x}_{-j}^{(i_2)}}{\sum_{i_2 \in I_{j,c_2}} w^{(i_2)}} \right\|_2^2,$$

where  $c_1, c_2$  are categorical variables from the set  $[K] := \{1, \dots, K\}$ . To apply CRLR on high-dimensional images, we adapt this categorical version on top of VQ-VAE discrete codes. The implementation details of the VQ-VAE are same as OREO (see Appendix B.2). Given a state-action pair  $(s^{(i)}, a^{(i)})$ , a VQ-VAE encoder  $g$  represents the state into code indices  $q^{(i)} := (q(i, 1), \dots, q(i, L))$  (see Section 3.1). The one-hot encoded version of the code indices  $q$  are denoted as  $\mathbf{q}$ , similarly to above. Then, a policy  $\pi$  and sample weights  $\{w^{(i)}\}$  are jointly trained by minimizing a weighted behavioral cloning objective and the proposed regularizer:

$$\begin{aligned} \mathcal{L}_{\text{CRLR}} &= \sum_i -w^{(i)} \log \pi(a^{(i)} | \mathbf{q}^{(i)}) \\ &+ \lambda \sum_j \sum_{\substack{c_1 < c_2 \\ \in \{c \in [K] \mid |I_{j,c}| > 0\}}} \left\| \frac{\sum_{i_1 \in I_{j,c_1}} w^{(i_1)} \mathbf{q}_{-j}^{(i_1)}}{\sum_{i_1 \in I_{j,c_1}} w^{(i_1)}} - \frac{\sum_{i_2 \in I_{j,c_2}} w^{(i_2)} \mathbf{q}_{-j}^{(i_2)}}{\sum_{i_2 \in I_{j,c_2}} w^{(i_2)}} \right\|_2^2, \end{aligned}$$

where  $\lambda$  is a loss weight for the regularizer. We update  $\pi$  and  $\{w^{(i)}\}$  iteratively until the objective converges, using the gradient descent optimizer.

## I Extended experimental results on confounded Atari environments

Table 9: Performance of policies trained on various confounded Atari environments without environment interaction. OREO achieves the best score on 15 out of 27 environments, and the best median and mean human-normalized score (HNS) over all environments. The results for each environment report the mean and standard deviation of returns over eight runs. CCIL<sup>†</sup> denotes the results without environment interaction.

Environment	BC	Dropout	DropBlock	Cutout	RandomShift	CCIL <sup>†</sup>	CRLR	OREO
Alien	954.1±83.9	1003.8±53.6	926.4±70.5	973.3±50.1	806.5±78.1	820.0±51.3	82.5±30.3	<b>1056.2±61.6</b>
Amidar	95.8±8.9	89.4±16.4	110.1±20.9	<b>118.7±14.0</b>	98.0±14.5	74.9±6.0	12.0±0.0	105.7±7.2
Assault	793.8±32.6	820.4±20.8	815.0±20.3	687.6±10.3	828.9±17.9	683.3±20.1	0.0±0.0	<b>840.9±27.8</b>
Asterix	292.2±168.4	313.8±115.3	345.4±207.7	212.4±83.2	135.5±85.4	64.3±2.8	<b>650.0±0.0</b>	180.8±65.2
BankHeist	442.1±20.7	485.7±19.7	508.4±14.2	486.1±14.2	367.2±17.7	<b>653.5±31.3</b>	0.0±0.0	493.9±17.6
BattleZone	11921.2±802.4	12457.5±427.7	12025.0±1425.9	11107.5±809.2	9180.0±592.3	6370.0±1227.5	1468.8±1512.6	<b>12700.0±1162.5</b>
Boxing	18.8±3.7	20.3±2.9	32.2±6.4	20.5±3.8	<b>38.3±5.2</b>	34.8±3.4	-43.0±0.0	36.4±5.0
Breakout	<b>5.7±0.5</b>	5.4±0.8	4.8±1.9	1.0±1.7	2.0±1.9	0.5±0.3	0.0±0.0	4.2±1.6
ChopperCommand	874.2±82.7	921.4±90.1	919.4±87.1	1016.1±169.0	936.4±125.6	760.6±58.7	<b>1077.2±9.1</b>	977.4±150.2
CrazyClimber	45372.9±5508.9	39501.6±6499.3	38345.6±7190.8	44523.2±8465.5	41924.0±7237.5	22616.8±3282.4	112.5±92.7	<b>55523.4±7722.2</b>
DemonAttack	157.2±12.5	180.5±21.5	167.8±12.2	173.1±11.3	<b>241.8±32.7</b>	171.3±17.3	0.0±0.0	224.5±45.4
Enduro	241.4±28.4	250.4±38.0	341.8±38.8	119.6±6.2	316.4±34.9	143.1±6.4	3.9±9.0	<b>522.8±29.1</b>
Freeway	32.3±0.1	32.4±0.2	32.7±0.1	32.5±0.2	33.0±0.1	<b>33.1±0.1</b>	21.4±0.1	32.7±0.2
Frostbite	116.3±21.1	124.5±26.7	128.2±35.6	<b>139.4±19.5</b>	121.6±16.4	53.3±30.7	80.0±0.0	129.9±12.8
Gopher	1713.9±182.5	1819.1±95.6	1818.2±150.3	1481.0±118.3	1995.0±189.1	1404.5±154.8	0.0±0.0	<b>2515.0±157.7</b>
Hero	11923.1±599.9	14109.7±894.1	14711.4±1119.9	14896.6±890.9	12816.0±988.2	6567.8±943.1	346.2±916.1	<b>15219.8±873.8</b>
Jamesbond	419.0±31.8	451.0±14.0	473.8±44.5	381.8±22.4	428.4±13.8	387.2±12.3	0.0±0.0	<b>502.8±39.3</b>
Kangaroo	2781.5±338.9	2912.9±266.5	3217.1±191.7	2824.0±200.8	1923.9±268.5	1670.5±153.4	122.8±215.4	<b>3700.2±126.0</b>
Krull	3634.3±70.6	3892.1±61.1	3832.1±281.0	3656.4±100.6	3788.7±216.3	3090.8±112.0	0.1±0.1	<b>4051.6±211.4</b>
KungFuMaster	15074.8±275.5	14452.1±865.4	15753.0±1265.2	11405.6±729.2	13389.9±624.3	13394.9±1261.9	0.0±0.0	<b>18065.6±1411.5</b>
MisPacman	1432.9±274.0	1733.1±273.2	1446.4±288.1	1711.0±184.6	1223.5±259.2	1084.2±199.1	105.3±60.5	<b>1898.4±229.8</b>
Pong	3.2±0.7	10.2±1.3	11.5±1.3	6.8±1.2	-0.1±2.2	-2.7±1.1	-21.0±0.0	<b>14.2±0.4</b>
PrivateEye	2681.8±270.2	2599.1±393.0	2720.6±427.4	2670.6±359.1	<b>3969.2±452.1</b>	305.3±247.5	-1000.0±0.0	3124.9±349.6
Qbert	5438.4±855.3	6469.0±760.3	6140.3±616.5	5748.6±655.5	3921.4±540.4	5138.0±437.9	125.0±0.0	<b>6966.4±443.5</b>
RoadRunner	18381.5±1519.9	21470.9±2274.4	22265.4±3168.3	12417.1±1307.8	16210.0±1193.1	11834.1±1936.3	1022.9±262.0	<b>24644.2±2235.1</b>
Seaquest	454.4±53.5	471.3±43.4	486.8±40.6	330.1±37.9	<b>1016.8±100.5</b>	271.2±11.5	172.5±19.8	753.1±63.6
UpNDown	4221.1±214.5	4147.1±426.2	<b>4789.2±201.0</b>	4159.6±585.5	3880.2±316.7	2631.1±224.0	20.0±0.0	4577.9±307.6
Median HNS	44.1%	47.4%	49.8%	42.0%	47.6%	36.2%	-1.5%	<b>51.2%</b>
Mean HNS	73.2%	79.0%	91.7%	69.5%	88.1%	71.7%	-45.9%	<b>105.6%</b>

## J Extended experimental results on original Atari environments

Table 10: Performance of policies trained on various original Atari environments without environment interaction. OREO achieves the best score on 14 out of 27 environments, and the best median and mean human-normalized score (HNS) over all environments. The results for each environment report the mean and standard deviation of returns over eight runs. CCIL<sup>†</sup> denotes the results without environment interaction.

Environment	BC	Dropout	DropBlock	Cutout	RandomShift	CCIL <sup>†</sup>	CRLR	OREO
Alien	986.5±54.4	1117.2±58.8	1094.8±73.7	1104.4±139.5	863.5±68.0	1050.4±62.4	100.0±0.0	<b>1222.2±95.4</b>
Amidar	90.8±7.7	81.6±8.2	113.5±12.9	125.0±7.7	78.2±9.2	78.6±3.1	12.0±0.0	<b>130.5±16.8</b>
Assault	816.8±25.0	901.1±22.6	829.9±23.7	694.1±9.5	848.7±17.0	755.5±9.9	0.0±0.0	<b>905.2±24.2</b>
Asterix	249.0±142.5	176.6±91.4	252.2±139.9	195.0±28.5	99.1±56.6	314.1±7.9	<b>592.5±148.4</b>	212.5±108.5
BankHeist	399.0±22.9	476.6±24.6	471.2±17.8	442.5±20.6	354.8±18.1	<b>606.1±31.7</b>	0.0±0.0	448.4±13.4
BattleZone	10933.8±642.0	11621.2±714.0	<b>12067.5±1269.0</b>	10641.2±328.5	8748.8±745.8	11191.2±709.5	5615.0±4482.6	11703.8±862.6
Boxing	21.8±4.6	25.7±4.1	32.1±5.0	21.2±3.4	35.8±4.3	34.2±2.9	-43.0±0.0	<b>39.9±2.2</b>
Breakout	<b>6.4±0.5</b>	2.9±2.5	6.0±0.9	3.1±2.4	4.4±2.4	2.1±2.0	0.0±0.0	5.4±1.0
ChopperCommand	1163.0±129.7	1162.0±51.9	1161.8±64.2	1183.9±56.4	1026.2±83.0	1027.2±78.2	1070.2±10.9	<b>1282.9±81.1</b>
CrazyClimber	54142.2±10143.4	54965.4±6305.6	55854.0±7056.0	47456.4±8129.0	60465.9±9050.9	39015.2±2266.3	885.5±864.8	<b>69380.1±8907.6</b>
DemonAttack	238.8±21.6	<b>359.3±47.3</b>	225.6±26.1	217.8±20.1	294.8±42.3	194.6±9.3	22.7±41.1	0.0±0.0
Enduro	226.2±24.6	304.6±31.4	359.1±38.0	132.9±4.7	282.2±27.4	182.8±6.2	0.8±1.2	<b>514.4±38.1</b>
Freeway	32.3±0.3	32.6±0.2	32.6±0.3	32.8±0.2	33.0±0.3	<b>33.1±0.2</b>	21.4±0.1	32.9±0.1
Frostbite	153.6±20.6	149.2±15.1	<b>165.7±19.7</b>	135.2±20.1	133.2±33.1	96.7±13.3	78.1±3.4	152.7±23.8
Gopher	1874.4±185.8	2220.4±156.2	2040.5±140.2	1588.2±106.1	1456.2±114.2	1301.9±219.5	0.0±0.0	<b>2903.9±146.6</b>
Hero	15100.4±774.6	15994.4±737.5	17058.6±419.4	15971.8±239.4	14867.2±904.5	<b>17487.6±813.5</b>	0.0±0.0	16370.3±501.4
Jamesbond	447.6±33.2	492.3±30.4	481.9±24.6	418.9±15.2	452.1±15.6	460.4±12.5	0.0±0.0	<b>527.9±20.7</b>
Kangaroo	3162.8±209.3	2860.4±175.1	<b>3638.6±312.6</b>	3242.6±124.2	2202.1±313.5	2938.1±391.6	0.0±0.0	3602.9±189.6
Krull	4447.9±91.5	<b>4764.7±112.3</b>	4526.5±113.7	4270.6±130.6	4611.6±144.9	4247.1±140.0	0.0±0.0	4633.6±114.9
KungFuMaster	12900.6±884.3	14994.5±1100.4	14819.0±806.0	9956.9±803.3	11698.0±1330.0	12876.9±912.2	0.0±0.0	<b>16955.5±1144.2</b>
MsPacman	1921.9±174.1	2022.6±202.8	2151.7±178.5	1949.7±176.1	1046.3±220.0	1160.6±144.1	70.0±0.0	<b>2263.8±165.3</b>
Pong	3.7±1.6	10.0±0.8	11.6±0.6	7.8±1.2	0.8±2.1	-19.8±0.4	-21.0±0.0	<b>12.5±0.7</b>
PrivateEye	3035.4±482.8	3396.3±205.9	3057.6±447.0	3092.2±305.9	<b>3578.9±222.9</b>	1016.4±286.8	-1000.0±0.0	3162.6±282.3
Qbert	5925.4±693.9	<b>6363.1±539.9</b>	5904.3±911.5	6174.8±585.8	4100.1±672.1	5056.3±456.9	125.0±0.0	5763.4±493.4
RoadRunner	18010.1±731.1	20137.8±1590.2	22522.5±1749.1	12698.9±1272.2	15615.4±712.1	18985.2±2105.5	1528.6±496.5	<b>27303.9±2326.7</b>
Seaquest	527.5±61.2	644.4±104.2	622.3±79.3	376.6±35.0	<b>948.0±95.5</b>	402.4±29.3	169.8±54.6	921.0±64.9
UpNDown	3782.1±245.7	3504.3±197.1	3886.4±257.1	3675.9±255.0	3500.4±246.8	3062.3±110.3	20.0±0.0	<b>4186.8±312.0</b>
Median HNS	46.7%	53.3%	47.7%	42.9%	47.3%	36.8%	-1.5%	<b>53.6%</b>
Mean HNS	82.0%	91.5%	99.0%	75.0%	91.7%	85.4%	-45.4%	<b>114.9%</b>