

# STUDYING THE INTERPLAY BETWEEN THE ACTOR AND CRITIC REPRESENTATIONS IN REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Extracting relevant information from a stream of high-dimensional observations is a central challenge for deep reinforcement learning agents. Actor-critic algorithms add further complexity to this challenge, as it is often unclear whether the same information will be relevant to both the actor and the critic. To this end, we here explore the principles that underlie effective representations for an actor and for a critic. We focus our study on understanding whether an actor and a critic will benefit from a decoupled, rather than shared, representation. Our primary finding is that when decoupled, the representations for the actor and critic systematically specialise in extracting different types of information from the environment—the actor’s representation tends to focus on action-relevant information, while the critic’s representation specialises in encoding value and dynamics information. Finally, we demonstrate how these insights help select representation learning objectives that play into the actor and critic’s respective knowledge specialisations, and improve performance in terms of agent returns.

## 1 INTRODUCTION

In recent years, auxiliary representation learning objectives have become increasingly prominent in deep reinforcement learning (RL) agents (Yarats et al., 2021; Dunion et al., 2024). These objectives facilitate extracting relevant features from high dimensional observations, and can help improve the sample efficiency and generalisation capabilities of both value-based (Anand et al., 2019; Schwarzer et al., 2021) and actor-critic methods (Yarats et al., 2019; Zhang et al., 2021; McInroe et al., 2023). However, knowing whether a particular representation learning objective will work and understanding *why* it works is often difficult due to the interplay between the different components of modern RL algorithms.

Online actor-critic algorithms like PPO (Schulman et al., 2017) jointly optimise policy improvement and value estimation objectives. When parametrised by deep neural networks, the actor (in charge of improving the policy) and the critic (in charge of estimating the value of the current policy) often share the same learned representation  $\phi$ , which maps observations to latent features  $z$ . While a coupled representation (Figure 1, top) reduces memory footprint and training costs, Cobbe et al. (2021) and Raileanu & Fergus (2021) report that fully decoupling the actor and critic (Figure 1, bottom) improves sample efficiency, and minimises overfitting to the environment instances, or *levels*, available during training.

In this study, we investigate three questions:

1. Why do decoupled representations achieve better performance?
2. Will the actor and critic benefit from specialised representation learning objectives?
3. What interplay remains between the actor and critic, once they are decoupled?

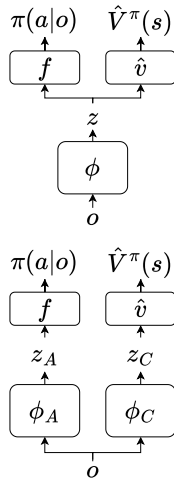


Figure 1: Models with shared (top) and decoupled representations (bottom).

Table 1: Once decoupled, the actor and critic representations  $\phi_A$  and  $\phi_C$  specialise in capturing different information from the environment. Reported values correspond to a PPO agent trained in Procgen (Cobbe et al., 2020). See §2 and §3 for formal definitions of the quantities quoted.

If ... is high,	it is possible to ...	% change from using a shared representation	
		$\phi_A$	$\phi_C$
$I(Z; L)$	overfit to training levels, i.e., use $z$ to identify levels.	-20%	+35%
$I(Z; V)$	use $z$ to predict state values.	+37%	+41%
$I((Z, Z'); A)$	use $z$ and $z'$ obtained from consecutive timesteps $t, t'$ to identify the action taken at timestep $t$ .	+23%	-48%
$I(Z; Z')$	differentiate between latent pairs obtained from consecutive and non-consecutive timesteps.	-96%	+324%

Our main findings are summarised below.

- Decoupled actor and critic representations extract different information about the environment. This information specialisation, described and quantified in Table 1, systematically occurs in the on-policy algorithms and benchmarks covered by our study, and is consistent with the actor’s and critic’s respective optimal representations.
- The actor benefits from representation learning when it prioritises extracting level-invariant information (i.e. features relevant to all levels) over level-specific information. This prioritisation matters more than picking a particular auxiliary objective.
- Through its role as a baseline in the actor’s objective, a decoupled critic will tend to bias policy updates to facilitate the optimisation of its own learning objective. The critic, therefore, plays an important role in exploration and data collection during training. Thus, we find that care must be taken when selecting a representation learning objective for the critic: certain objectives improve the critic’s value predictions but may prevent convergence to the optimal policy because the objective induces significant bias.

## 2 BACKGROUND

**Zero-shot transfer in RL.** We consider the problem of zero-shot transfer in RL in the episodic setting. Following the framework established by Kirk et al. (2023), we model the environment as a Contextual-MDP (CMDP)  $\mathcal{M} = (\mathbb{S}, \mathbb{A}, \mathbb{O}, \mathcal{T}, \Omega, R, \mathbb{C}, P(c), \mathcal{P}_0, \gamma)$  with state, action and observation spaces  $\mathbb{S}$ ,  $\mathbb{A}$  and  $\mathbb{O}$  and discount factor  $\gamma$ . In a CMDP, the reward  $R : \mathbb{S} \times \mathbb{C} \times \mathbb{A} \rightarrow \mathbb{R}$ , and observation functions  $\Omega : \mathbb{S} \times \mathbb{C} \rightarrow \mathbb{O}$  as well as the transition  $\mathcal{T} : \mathbb{S} \times \mathbb{C} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{S})$ , and initial state  $\mathcal{P}_0 : \mathbb{C} \rightarrow \mathcal{P}(\mathbb{S})$  kernels can change with the *context*  $c \in \mathbb{C}$ , with  $c \sim P(c)$  at the start of each episode. The CMDP is therefore conceptually equivalent to an MDP with state space  $\mathbb{X} : \mathbb{S} \times \mathbb{C}$ . Each context  $c$  maps one-to-one to a particular environment instance, or *level*, and thus represents the component of the state  $x$  that cannot change during the episode. During training, we assume access to a limited set of training levels  $L$ , and we measure transfer by evaluating the RL agent on an held-out set  $L_{\text{test}}$  (both obtained by sampling from  $P(c)$ ). The agent’s policy  $\pi : \mathbb{O} \rightarrow \mathcal{P}(\mathbb{A})$  maps observations to action distributions and induce a value function  $V^\pi : \mathbb{X} \rightarrow \mathbb{R}$  mapping states to expected future returns  $V^\pi(x) = \mathbb{E}_\pi[\sum_t^T \gamma^t r_t]$ , where  $\{r_t\}_{0:T}$  are possible sequences of rewards obtainable when following policy  $\pi$  from  $x$  and until the episode terminates. We define the optimal policy for zero-shot transfer  $\pi^*$  as the policy maximising  $\mathbb{E}_{c \sim P(c), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)]$ .

**Actor-critic architectures.** On-policy actor-critic models consist of a policy network  $\pi_{\theta_A}$  and a value network  $\hat{V}_{\theta_C}$ , with *actor* parameters  $\theta_A$  and *critic* parameters  $\theta_C$  (we use  $\cdot_A/\cdot_C$  when referring to the actor/critic in this work). When learning from high dimensional observations, such as pixels, a representation  $\phi : \mathbb{O} \rightarrow \mathbb{Z}$  maps observations to latent features  $z \in \mathbb{Z}$ . When coupled, the policy and value networks share a representation and split into actor and critic heads  $f$  and  $\hat{v}$ . That is, we

have  $\pi_{\theta_A} = f_{\omega} \circ \phi_{\eta}$  and  $\hat{V}_{\theta_C} = \hat{v}_{\xi} \circ \phi_{\eta}$ , with  $\theta_A = (\omega, \eta)$  and  $\theta_C = (\xi, \eta)$ . When decoupled, two representation functions  $\phi_A, \phi_C$  with parameters  $(\eta_A, \eta_C)$  are learned.

**PPO and PPG.** In this work, we investigate the representation properties of PPO (Schulman et al., 2017) and Phasic Policy Gradient (PPG) (Cobbe et al., 2021), two actor-critic algorithms that have been reported to benefit from improved sample efficiency and transfer upon decoupling (Raileanu & Fergus, 2021; Cobbe et al., 2021). In PPO, the actor maximises

$$J_{\pi}(\theta_A) = \mathbb{E}_B \left[ \min \left( \frac{\pi_{\theta_A}(a_t|o_t)}{\pi_{\theta_{A_{old}}}(a_t|o_t)} \hat{A}_t, \text{clip} \left( \frac{\pi_{\theta_A}(a_t|o_t)}{\pi_{\theta_{A_{old}}}(a_t|o_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) + \beta_H \mathbf{H}(\pi_{\theta_A}(a_t|o_t)) \right], \quad (1)$$

where  $\theta_{A_{old}}$  are the actor weights before starting a round of policy updates,  $B$  is a batch of trajectories collected with  $\pi_{\theta_{A_{old}}}$ ,  $\hat{A}_t$  is an estimator for the advantage function at timestep  $t$ ,  $\mathbf{H}(\cdot)$  denotes the entropy and  $\epsilon$  and  $\beta_H$  are hyperparameters controlling clipping and the entropy bonus. The critic minimises

$$\ell_V(\theta_C) = \frac{1}{|B|} \sum_{o_t \in B} (\hat{V}_{\theta_C}(o_t) - \hat{V}_t)^2, \quad (2)$$

where  $\hat{V}_t$  are value targets. Both  $\hat{A}$  and  $\hat{V}$  are computed using GAE (Schulman et al., 2016). PPG performs an auxiliary phase after conducting PPO updates over  $N_{\pi}$  policy phases. To prevent overfitting, the auxiliary phase fine-tunes the critic and distills value information into the representation from much larger trajectory batches  $B_{aux} = \bigcup_{i \in 1, \dots, N_{\pi}} B_i$ , using the loss  $\ell_{\text{joint}} = \ell_V + \ell_{aux}$ , with

$$\ell_{aux}(\theta_A) = \frac{1}{|B_{aux}|} \sum_{(a_t, o_t) \in B_{aux}} (\hat{V}_{\theta_A}^{aux}(o_t) - \hat{V}_t)^2 + \beta_c D_{\text{KL}}(\pi_{\theta_{A_{old}}}(a_t|o_t) \| \pi_{\theta_A}(a_t|o_t)), \quad (3)$$

where  $\beta_c$  controls the distortion of the policy. When decoupled,  $\hat{V}^{aux} = v^{aux} \circ \phi_A$  distills value information into representation parameters  $\eta_A$  through an additional head  $v^{aux}$ . When coupled,  $v^{aux} = \hat{v}$ , and a stop-gradient operation on  $\ell_V$  ensures  $\eta$  is updated by the critic during the auxiliary phase only.

**Mutual information.** We study the information embedded in features  $z$  outputted by  $\phi$ . To do so, we propose metrics based on the mutual information  $\mathbf{I}(X; Y)$ , measuring the information shared between sets of random variables  $X$  and  $Y$ , defined as

$$\mathbf{I}(X; Y) = \mathbf{H}(X) + \mathbf{H}(Y) - \mathbf{H}(X, Y) = \sum_{\mathbf{X}} \sum_{\mathbf{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}, \quad (4)$$

where integrals replace sums for continuous quantities.  $\mathbf{I}(X; Y)$  is symmetric, and quantifies how much information about  $Y$  is obtained by observing  $X$ , and vice versa. Similarly, the conditional mutual information  $\mathbf{I}(X; Y|Z)$  measures the information shared between  $X$  and  $Y$  that does not depend on  $Z$ . Finally, our results build upon the data-processing inequality, which states that when  $X, Y$  and  $Z$  form the Markov chain  $X \rightarrow Y \rightarrow Z$  we have  $\mathbf{I}(X; Z) \leq \mathbf{I}(X; Y)$ . Simply put, all information that “flows” from  $X$  to  $Z$  must first flow through  $Y$ .

We measure mutual information using the k-nearest neighbors entropy estimator proposed by Kraskov et al. (2004) and extended to pairings of continuous and discrete variables by Ross (2014). We briefly introduce notation for random variables used in following sections.  $L \sim P(c)$  denotes the set of training levels drawn from the CMDP context distribution.  $A, R, O, O', X$  and  $X'$  are sets constructed from  $n$  transitions  $(a_t, r_t, o_t, o_{t+1}, x_t, x_{t+1})$  uniformly sampled from a batch of trajectories collected in  $L$  using policy  $\pi$ .  $Z$  and  $Z'$  are latent features, with  $z = \phi(o), z' = \phi(o')$ . We construct  $V$  using the rewards obtained from  $t$  until episode termination, with  $v_t = \sum_{\bar{t}=t}^T \gamma^{\bar{t}-t} r_{\bar{t}}$ .

### 3 CATEGORISING AND QUANTIFYING THE INFORMATION EXTRACTED BY LEARNED REPRESENTATIONS

To conduct our analysis of the respective functions of the actor and critic representations, we analyse the information being extracted from observations at agent convergence. We propose four mutual information metrics to measure information extracted about the identity of the current training level,

the value function and the inverse and forward transition dynamics of the environment. This categorisation will not completely capture what information gets extracted by representations. However, we discover it still highlights a specialisation for the actor’s and critic’s representations. Moreover, our proposed categorisation is strengthened by different theoretical arguments, which are discussed next.

**Overfitting.** Our first metric,  $I(Z; L)$ , quantifies overfitting of the actor and critic representations to the set of training levels, as it measures how easy it is to infer the identity of the current level from  $Z$ . We follow the same reasoning as Garcin et al. (2024) to derive an upper bound for the generalisation error that is proportional to  $I(Z_A; L)$ .<sup>1</sup>

**Theorem 3.1.** *The difference in returns achieved in train levels and under the full distribution, or generalisation error, has an upper bound that depends on  $I(Z_A; L)$ , with*

$$\mathbb{E}_{c \sim \mathcal{U}(L), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] - \mathbb{E}_{c \sim P(c), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] \leq \sqrt{\frac{2D^2}{|L|}} \times I(Z_A; L), \quad (5)$$

where  $c \sim \mathcal{U}(L)$  indicates  $c$  is sampled uniformly over levels in  $L$ ,  $D$  is a constant such that  $|V^\pi(x)| \leq D/2, \forall x, \pi$  and  $Z_A$  is the output space of the actor’s learned representation.

**Value information.** The second metric quantifies  $I(Z; V)$ , the mutual information between  $Z$  and the state values in  $L$  when following  $\pi$ . A high  $I(Z_C; V)$  should help optimise  $\ell_V$  (Equation (2)). However, we wish to understand whether increasing  $I(Z_A; V)$  is always desirable: Cobbe et al. (2021) report that value distillation into the actor’s representation improves sample efficiency over  $L$ , whereas Raileanu & Fergus (2021) and Garcin et al. (2024) report that a coupled PPO agent achieves stronger generalisation over  $L_{\text{test}}$  when  $\ell_V$  remains high during training.

**Dynamics in the latent space.** The remaining two metrics investigate the transition dynamics  $\mathcal{T}_z : \mathbb{Z} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{Z})$  within the latent state space  $\mathbb{Z}$  spanned by the representation. We will see in later sections that the *reduced MDP*  $(\mathbb{Z}, \mathbb{A}, \mathcal{T}_z, R_z, \gamma)$  spanned by the actor or critic’s representation tend to have distinct  $\mathcal{T}_z$ , which often markedly differ from the transition dynamics  $\mathcal{T}$  in the original environment.  $I(Z; Z')$  measures whether it is possible to differentiate latent-pairs  $(\phi(o), \phi(o'))$  obtained from consecutive observations from latent-pairs from non-consecutive observations.  $I((Z, Z'); A)$  quantifies how easy it is to predict the action given the pair  $(\phi(o), \phi(o'))$ . In Theorem 3.2, we establish that  $\mathcal{T}_z$  maintains the *Markov property* of the original MDP when both of these metrics attain their theoretical maximum.<sup>2</sup>

**Theorem 3.2.** *if  $\mathcal{T} : \mathbb{X} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{X})$  satisfies the Markov property, and we have  $I((X, X'); A) = I((Z, Z'); A)$  and  $I(X; X') = I(Z; Z')$  for any  $X, X', A, Z, Z'$  collected using policy  $\pi$ , then  $\mathcal{T}_z : \mathbb{Z} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{Z})$  satisfies the Markov property when following  $\pi$ .  $\mathcal{T}_z$  always satisfies the Markov property if the above conditions hold for any  $\pi$ .*

Given that  $\phi$  only induces  $\mathcal{T}_z$  for the current  $\pi$  in the on-policy setting, we make the distinction between  $\mathcal{T}_z$  being Markov when following  $\pi$  and the more general notion of  $\mathcal{T}_z$  being Markov when following any policy. Crucially, Theorem 3.2 generalises the equivalence relations obtained by Allen et al. (2021) to continuous metrics.<sup>3</sup> As such,  $I((Z, Z'); A)$  and  $I(Z; Z')$  quantify how close any  $\phi$  comes to have  $\mathcal{T}_z$  satisfy the Markov property. They remain useful in settings in which it isn’t practical (or even possible) for  $\mathcal{T}_z$  to satisfy the Markov property (e.g. when observations are high dimensional, and/or under partial observability).

## 4 INFORMATION SPECIALISATION IN ACTOR AND CRITIC REPRESENTATIONS

Raileanu & Fergus (2021); Cobbe et al. (2021) have attributed the performance improvements obtained from decoupled architectures to the disappearance of gradient interference between the actor

<sup>1</sup>Proofs for the theoretical results presented in this work are provided in Appendix A.

<sup>2</sup>The Markov property is satisfied for a MDP  $(\mathbb{Z}, \mathbb{A}, \mathcal{T}_z, R_z, \gamma)$  if and only if  $\mathcal{T}_z^{(k)}(z_{t+1} | \{a_{t-i}, z_{t-i}\}_{i=0}^k) = \mathcal{T}_z(z_{t+1} | a_t, z_t)$  and  $R_z^{(k)}(z_{t+1} | \{a_{t-i}, z_{t-i}\}_{i=0}^k) = R_z(z_{t+1} | a_t, z_t), \forall a \in \mathbb{A}, z \in \mathbb{Z}, k \geq 1$ .

<sup>3</sup>Allen et al. (2021) consider Block-MDPs (Du et al., 2019), which are MDPs in which observations are guaranteed to maintain the Markov property.

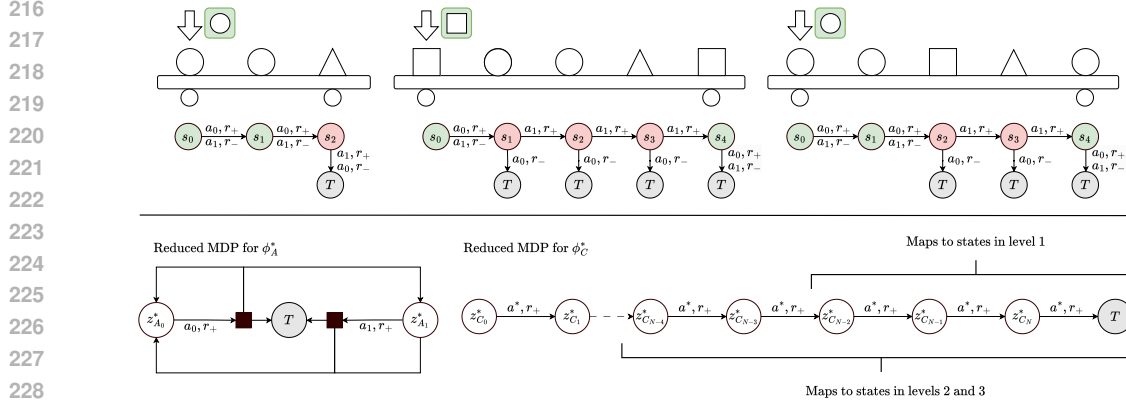


Figure 2: (Top) the initial observations and state spaces of three levels from the assembly line environment in §4.1. (Bottom) the reduced MDPs spanned by  $\phi_A^*$  and  $\phi_C^*$ .

and critic, and to the critic tolerating a higher degree of sample reuse than the actor before overfitting. We propose a different interpretation: given their different learning objectives, the actor’s and critic’s *optimal representations* (defined below) prioritise different types of information from the environment. While not incompatible with prior interpretations, our claim is stronger. We posit that an optimal (or near-optimal) representation for both the actor and critic will generally be impossible under a shared architecture.

**Definition 4.1.** Given the model  $m = f_\omega \circ \phi$  and associated loss  $\ell_m(\omega, \phi)$ , an optimal representation  $\phi^* : \mathbb{O} \rightarrow \mathbb{Z}^*$  satisfies the conditions:

1. **Optimality conservation.**  $\min_\omega \ell_m(\omega, \phi^*) = \min_{\omega, \phi} \ell_m(\omega, \phi)$
2. **Maximal compression.**  $\phi^* \in \arg \min_{\tilde{\Phi}} |\mathbb{Z}^*|$ , with  $\tilde{\Phi}$  the set of all  $\phi$  satisfying condition 1.

#### 4.1 WARMUP: AN ASSEMBLY LINE INSPECTION PROBLEM

Here we present a motivating example to highlight the respective specialisations and mutual incompatibility of  $\phi_A^*$  and  $\phi_C^*$ . Starting from this example, we derive several conditions for specialisation that apply irrespective of the setting.

In our example, the agent inspects parts for defects on an assembly line. The agent is trained on a set  $L$  of levels drawn from  $P(c)$ . A level is characterised by a particular combination of part specifications, number and ordering, each part having a probability  $P^F$  of being defective. We depict three possible levels in Figure 2. At each timestep, the agent observes the part specifications for the current level, which parts are on the assembly line and which part is up for inspection. The agent picks action  $a \in \mathbb{A} = \{a_0 = \text{accept}, a_1 = \text{reject}\}$  and moves to the next part. It receives a reward  $R = r_+$  when correctly accepting/rejecting a good/defective part and  $R = r_-$  when it makes a mistake, with  $r_+ > r_-$ . The episode terminates early when the agent accepts a defective part, otherwise it terminates after  $N_c$  timesteps, where  $N_c$  is the number of parts in level  $c \in L$ .

##### 4.1.1 THE ACTOR’S OPTIMAL REPRESENTATION

The combinatorial explosion of possible specifications and part assortments means  $\phi_A^*$  should ideally map observations to a *reduced MDP* spanning a much smaller state space than in the original environment. However,  $\phi_A^*$  should still provide the information necessary to select the optimal action at each timestep of each level, including those not in the training set.

**Dynamics of the reduced MDP.** Under our definition, the mapping

$$\phi_A^*(o) = \begin{cases} z_{A0}^*, & \text{if } a^* = a_0 \\ z_{A1}^*, & \text{if } a^* = a_1. \end{cases} \quad (6)$$

spans the reduced MDP in Figure 2 (bottom left), which describes the perceived environment dynamics when only observing the latent states in  $Z_A^*$ . By construction,  $I((Z_A^*, Z_A^{*'}); A)$  is guaranteed to be maximised when following the optimal policy. On the other-hand, encoding information about forward dynamics is not necessary for optimality: even if transitions are deterministic and  $I(X; X')$  is maximised under  $\pi^*$  in the original environment, in this reduced MDP we have  $I(Z_A^*; Z_A^{*'}) = 0$ .

**Overfitting to training levels.**  $I(Z_A^*; L)$  should be small for  $\phi_A^*$  to be invariant to individual levels and allow zero-shot transfer, and this is made evident in our example. However, we can picture an *overfit*  $\phi_A$  with high  $I(Z_A; L)$ , that encodes spurious correlations, to first identify, and then solve certain levels. For example: “if the rightmost object is a triangle, then I must be in level 1, and, in level 1, only the triangle has a defect”.

It is therefore desirable to minimise  $I(Z_A^*; L)$ . However, we note that even  $\phi_A^*$  does not perfectly achieve this, as it satisfies sufficient conditions for which  $I(Z_A^*; L)$  must be positive. We provide these conditions in Lemma 4.1.

**Lemma 4.1.**  $I(Z; L) > 0$  if  $\exists z_k, c_j \in Z \times L$  such that  $\mu(z_k|c_j) \neq \mu(z_k)$  and  $I(O; L) > 0$ ,  $I(O; L) > 0$  being the mutual information between  $L$  and observations  $O$ , with  $\phi(o) = z \in Z$ .

In our example,  $I(O; L) > 0$ , since each observation depicts quantities unique to each level. We can confirm the first condition by inspecting the stationary distributions in a particular level  $c$  and over all levels,

$$\mu(z) = \begin{cases} \bar{P}^F, & \text{if } z = z_{A0}^* \\ P^F, & \text{if } z = z_{A1}^* \end{cases} \quad \mu(z|c) = \begin{cases} \bar{P}_c^F, & \text{if } z = z_{A0}^* \\ P_c^F, & \text{if } z = z_{A1}^* \end{cases}, \quad (7)$$

where  $\bar{P}^F = 1 - P^F$  and  $P_c^F$  is the defect probability when in level  $c$ . While  $\mathbb{E}_c[P_c^F] = P^F$ , individual levels will not all have the same distribution of defective parts. For example, being in  $z_{A1}^*$  makes it more likely to be in the second out of the three levels depicted in Figure 2, since it is where  $\mu(z_{A1}^*)$  is the highest.

A representation  $\phi_A$  with a high  $I(Z_A; Z'_A)$  may also indirectly encode information relevant to specific levels. We show in Lemma 4.2 that  $I(Z; L)$  increases when  $I(Z; Z')$  increases, whenever the information gained does not apply to all levels in  $L$ .

**Lemma 4.2.**  $I(Z; L)$  monotonically increases with  $I(Z; Z') - I(Z; Z'|L)$ .

For example,  $\phi_A$  encoding that “An object that comes after two spheres always has a defect” applies to the three levels in Figure 2, but not to all possible levels, and  $I(Z; L)$  increases because  $I(Z; Z') > I(Z; Z'|L)$ . Conversely,  $\phi_A$  encoding that “the arrow above an object moves to the right each timestep” would not increase  $I(Z; L)$ , because, as this information applies to all levels, we have  $I(Z; Z') = I(Z; Z'|L)$ .

The key implications of these lemmas and the above examples are that 1) When  $I(Z_A; L)$  is high, it is possible for  $\phi_A$  to be optimal over  $L$ , but not over unseen levels, 2)  $\phi_A^*$  being optimal does not always guarantee  $I(Z_A^*; L) = 0$  and zero-generalisation error 3) High  $I(Z_A; Z'_A)$  may cause overfitting, due to its relationship with  $I(Z_A; L)$ .

**Value distillation may induce overfitting.** We now investigate the information  $\phi_A^*$  encodes about the value function. Lemma 4.3 establishes a sufficient condition for  $I(Z; V)$  to be positive. We then use this condition in Corollary 4.4 to show that  $Z_A^*$  being invariant to rewards by construction isn’t sufficient for  $\phi_A^*$  to not extract any information about the value function.

**Lemma 4.3.**  $I(Z; V) > 0$  if  $\exists z_k, v_m \in Z \times V$  such that  $\frac{1}{L} \sum_{c \in L} p(z_k, v_m|c) \neq p(z_k)p(v_m)$ .

**Corollary 4.4.**  $I(Z; V)$  can be positive when  $Z|L$  and  $V|L$  are conditionally independent. If  $I(Z; V) > 0$  and  $Z|L$  and  $V|L$  are conditionally independent, then  $I(Z; L) > 0$ .

Figure 2 showcases an example of this. States values can be higher in levels 2 and 3 than in level 1 because they have more timesteps. Those levels have a higher chance of being in  $z_{A1}^*$  because  $P_c^F$  is higher. Then, both the state values and optimal action distributions share a confounder in the level identities, and lead to  $I(Z_A^*; V)$  being positive. This could challenge the notion that value distillation is an effective way to improve the actor’s representation, as 1) value information is not necessary for  $\phi_A^*$  to be optimal and 2) to minimise the value distillation loss, the agent may increase  $I(Z; V)$  by using  $I(Z; L)$  as a proxy, i.e. by overfitting to  $L$ .

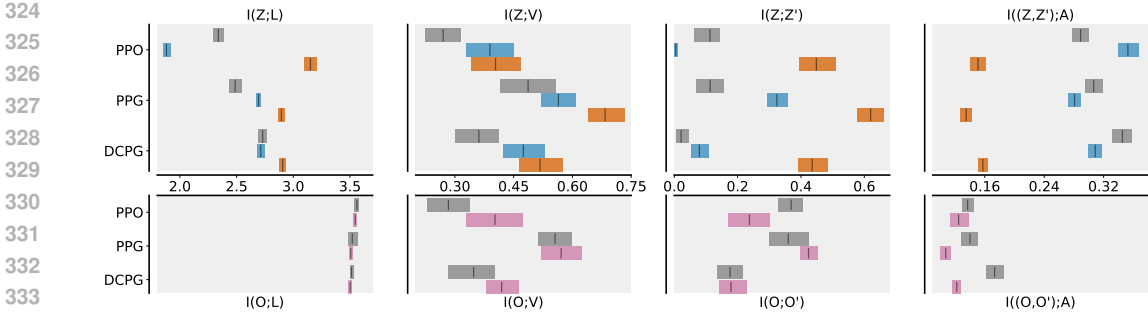


Figure 3: Mean and 95% confidence interval of  $I(Z; \cdot)/I(O; \cdot)$  (top/bottom) in Procgen. Top row shows shared (gray), actor (blue), and critic (orange) representations. Bottom row shows shared (gray) and decoupled (pink). X-axes are shared across top and bottom. For all algorithms, decoupling induces specialisation consistent with §4.

#### 4.2 THE CRITIC’S OPTIMAL REPRESENTATION

The reduced MDP spanned by  $\phi_C^*$  is depicted in Figure 2 (bottom right). In order to ensure perfect value prediction,  $\phi_C^*$  maps each possible optimal state value to a different element in  $\mathbb{Z}_C^*$ , and it maximises  $I(\mathbb{Z}_C^*; V)$  by construction.  $I(\mathbb{Z}_C^*; Z_C^{*'})$  is also high due to the recurrence  $V^\pi(x) = \mathbb{E}_{a \sim \pi}[R(x, a) + \gamma \mathbb{E}_{x' \sim P^\pi(x'|x)}[V^\pi(x')]]$ . Lemma 4.2 tells us that  $V^\pi$  is a quantity inherently more level specific than the optimal action for one timestep, because  $V^\pi$  encodes information pertaining to all future timesteps. Therefore, we should expect that, in general,  $I(\mathbb{Z}_C^*; L) > I(\mathbb{Z}_A^*; L)$ . Corollary 4.5 tells us that  $I((\mathbb{Z}_C^*, Z_C^{*'}); A)$  is similar to  $I(\mathbb{Z}_A^*; V)$  in the sense that, while  $\mathbb{Z}_C^*$  should be invariant to actions, we may still observe positive  $I((\mathbb{Z}_C^*, Z_C^{*'}); A)$  due to confounding induced when  $I(\mathbb{Z}_C^*; L) > 0$ .

**Corollary 4.5.**  $I((Z; Z'); A)$  can still be positive when  $(Z, Z')|L$  and  $A|L$  are conditionally independent. If  $I((Z; Z'); A) > 0$  and  $(Z, Z')|L, A|L$  are conditionally independent, then  $\{I(Z; L) > 0$  and  $I(A; L) > 0\}$  and/or  $\{I(Z'; L) > 0$  and  $I(A; L) > 0\}$ .

$\phi_C^*$  is not compatible with  $\pi^*$ . Paradoxically, while  $\phi_C^*$  would necessitate trajectories collected using the optimal policy in order to be learnt, it is not possible to have an optimal policy that only depends on  $z_C^*$ . We trace this issue back to  $\phi_C^*$  not satisfying the first condition of Theorem 3.2 under  $\pi^*$ . The information contained in  $z_C^*$  is not sufficient for picking the optimal action in any given timestep, and therefore the best response is to always pick  $a_1$  in order to prevent early termination.

#### 4.3 CONFIRMING SPECIALISATION IN THE PROCGEN BENCHMARK

We conclude this section by studying the representations learned by PPO (Schulman et al., 2017), PPG (Cobbe et al., 2021) and DCPG (Moon et al., 2022), a close variant of PPG that employs delayed value targets to train the critic and for value distillation. We evaluate all algorithms with and without decoupling their representation. We conduct our experiments in Procgen (Cobbe et al., 2020), a benchmark of 16 games designed to measure generalisation in RL. We report our main observations in below, with extended results and details on our methodology included in Appendix C.2.

**Specialisation is consistent with  $\phi_A^*$  and  $\phi_C^*$ .** As no algorithm achieves optimal scores in all games, we now consider the suboptimal representations  $\phi_A$  and  $\phi_C$  realistically obtainable by the end of training. In Figure 3, we observe clear specialization upon decoupling consistent with the properties we expect for  $\phi_A^*$  and  $\phi_C^*$ .  $\phi_C$  has high  $I(Z; V)$ ,  $I(Z; Z')$  and  $I(Z; L)$ , while  $\phi_A$  specializes in  $I((Z, Z'); A)$ .

**Decoupling is more parameter efficient.** Since decoupled representations fit twice as many parameters, it is fair to wonder whether the performance improvements are caused by an increased model capacity. To test this, we measure performance as we scale model size in a shared and a decoupled architecture in Figure 4. Surprisingly, the decoupled agent outperforms a shared model with four times its original parameter count. This makes the decoupled architecture at least twice as parameter efficient as a shared architecture.

**On Markov representations.** Theorem 3.2 tells us that a representation is Markov when  $I((Z, Z'); A)$  and  $I(Z; Z')$  are both maximised. Yet, sometimes upon decoupling,  $I((Z_A, Z'_A); A)$  increases and  $I(Z_A; Z'_A)$  decreases, while  $I((Z_C, Z'_C); A)$  decreases and  $I(Z_C; Z'_C)$  increases. This suggests that neither the actor nor the critic particularly benefit from a Markov representation. It is consistent with the fact that neither  $\phi_A^*$  nor  $\phi_C^*$  need to be Markov to be optimal. In fact, we find no clear correlation between  $I((Z, Z'); A) + I(Z; Z')$  (Figure 7) and agent performance (Figure 13).

## 5 REPRESENTATION LEARNING FOR THE ACTOR

In this section, we study how different representation learning objectives affect  $\phi_A$  in PPO, PPG and DCPG. We consider advantage (Raileanu & Fergus, 2021) and dynamics (Moon et al., 2022) prediction, data augmentation (Raileanu et al., 2021) and MICo (Castro et al., 2021), an objective explicitly shaping the latent space to embed differences in state values. We study these objectives in Procgen (Figure 5), and in four continuous control environments with video distractors, similar to those from Stone et al. (2021), which we re-implement in Brax (Freeman et al., 2021) (Figure 11).

**Representation learning impacts information specialisation.** As expected, applying auxiliary tasks alters what information is extracted by the representation. With few exceptions, dynamics prediction plays into the information specialisation of  $\phi_A$  by consistently increasing  $I((Z, Z'); A)$  and reducing  $I(Z; Z')$ ,  $I(Z; V)$ , and  $I(Z; L)$ . On the other-hand, MICo has the opposite effect on the aforementioned quantities (sans  $I(Z; V)$  in Procgen and PPO in Brax for  $I((Z, Z'); A)$ ). The effects of the last two objectives are not as clear-cut. Data augmentation produces little change in each quantity, while advantage prediction tends to reduce the measured mutual information, but is inconsistent in the quantities it affects. Performance-wise, data augmentation improves train and test scores for all algorithms; dynamics prediction tends to improve performance for PPG and DCPG; MICo tends to decrease performance, and advantage prediction makes no noticeable impact. We hypothesise that an effective representation learning objective for  $\phi_A$  does not change its specialisation (data augmentation) or plays into it (dynamics prediction), and does not work against it (MICo).

**On the importance of the batch size and data diversity.** We now turn our attention to an apparent contradiction in the relationship between value distillation and performance. Decoupling PPO, and thus completely forgoing value distillation, leads to improved train and test scores (Figure 13). However, PPG and DCPG conduct four times as many value distillation updates as coupled PPO during training, and achieve an even more significant performance improvement. Crucially, conducting value distillation every  $N_\pi$  policy phases ensures the batch size  $B_{\text{aux}}$  is  $N_\pi$  times as large as the batch size used in PPO (32 times in our experiments), which increases data diversity. We hypothesise that this diversity is key for successful representation learning: it leads to more information being embedded in  $\phi_A$  and promotes *learning invariances* to seemingly unrelated quantities, improving generalisation. To test this hypothesis, we train PPG under different  $B_{\text{aux}}$  while keeping  $N_\pi$  and the total number of distillation updates the same. We report scores,  $I(Z_A; V)$ , and  $I(Z_A; L)$  in Figure 12. The agent’s scores are proportional to  $B_{\text{aux}}$ , consistent with a similar experiment conducted by Wang et al. (2023). By the chain rule of mutual information, we can decompose  $I(Z_A; V)$ :

$$I(Z_A; V) = I(Z_A; V|L) + I(Z_A; L) - I(Z_A; L|V), \quad (8)$$

where  $I(Z_A; V)$  is the level-invariant information  $Z_A$  encodes about  $V$ , and  $I(Z; L) - I(Z; L|V)$  is the level-specific information encoded about  $V$ . We see that  $I(Z_A; V)$  keeps increasing with the batch size, while  $I(Z_A; L)$  plateaus. Put together, these results point to level-invariant information being prioritised when training on larger and more diverse batches of data, and this prioritisation being important for improving performance.

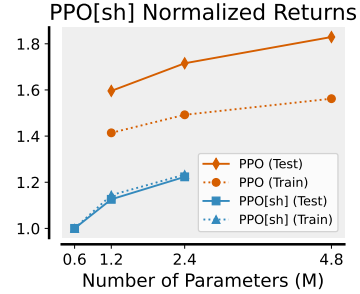


Figure 4: Parameter scaling experiments between coupled (blue) and decoupled (orange) PPO in Procgen.



## 6 THE CRITIC’S OBJECTIVE(S) WILL INFLUENCE DATA COLLECTION

We now consider how the same set of representation learning objectives affect the critic’s representation and present our results in Figures 8 and 11. The effect of a given objective on the information extracted by  $\phi_C$  is consistent with how they would have affected  $\phi_A$  in the previous section. However, we report two surprising findings: 1) Without conducting any value distillation, decoupled PPO has a 37% higher  $I(Z_A; V)$  than shared PPO (Table 1), and 2) the information specialisation of  $\phi_C$  incurred by applying an objective on the critic is often observed in  $\phi_A$ , albeit to a lesser extent. Given that the two representations are decoupled, how can an objective applied to  $\phi_C$  affect  $\phi_A$ ?

As we maintain different optimisers for the actor and critic, their only remaining interaction in decoupled PPO is through  $J_\pi$  (Equation (1)):  $\hat{A}_t$  being computed from the critic’s value estimates. Therefore, at least one of the following hypothesis must hold:

- Data collection bias.** Through  $J^\pi$  updates, the critic biases  $\pi$  to collect trajectories containing information relevant to its own learning objective. This information could then leak through  $\phi_A$  because more of this information is contained in its input. In this scenario, it is not necessary for  $\phi_A$  to become more proficient at extracting critic-relevant information.
- Implicit knowledge transfer.** The advantage targets in  $J^\pi$  induce information transfer between  $\phi_C$  and  $\phi_A$  when applying the gradients  $\nabla_{\theta_A} J^\pi$ . Here,  $\phi_A$  becomes proficient at extracting the same information  $\phi_C$  extracts.

The first hypothesis broadly holds in our experiments: in most cases, applying MICO to the critic increases  $I(O; V)$  and  $I(O; O')$ , and applying dynamics prediction increases  $I((O, O'); A)^4$ . Furthermore,  $I(O; V)$  increases when PPO is decoupled (Figure 3). Without the critic’s influence, there would be no direct incentive for the actor to collect data that contains value information, since no value distillation is taking place. It hints at interesting ramifications: the critic may have a much larger impact on exploration than previously imagined, and it may play fundamentally different roles between the offline and online RL settings. We leave the exploration of these ramifications to future work.

To test the second hypothesis, we measure the *compression efficiency*, applicable whenever  $I(O; \cdot) > 0$ , and defined as

$$C(Z; \cdot) = \min \left( \frac{I(Z; \cdot)}{I(O; \cdot)}, 1 \right). \quad (9)$$

For example,  $C(Z_A; V)$  measures the fraction of available information in  $I(O; V)$  that is extracted by  $\phi_A$ .<sup>5</sup> In Tables 2 and 3, we report that  $C(Z_A; V)$  does not significantly change between PPO[sh] and decoupled PPO, or when MICO is applied to the critic in decoupled PPO. This result appears to disprove the second hypothesis, at least in PPO. We cannot confirm whether implicit knowledge transfer occurs PPG and DCPG, as explicit knowledge transfer from the critic into the actor already occurs through value distillation.

Finally, we highlight that this data collection bias generally leads to worse performance (Figure 13), even with a representation learning objective aligned with the critic’s specialisation. Interestingly, employing different representation learning objectives for the actor and the critic results in surprising interactions. In Procgen, using advantage distillation on the actor does not affect performance, and using MICO on the critic degrades it. However, combining the two brings  $I(O; V)$  (Figure 10) back down to normal levels, and sharply increases PPO’s performance on both the train and test sets (Figure 13), suggesting that objectives aligned with  $\phi_A$  can reduce the bias in data collection induced by the critic.

<sup>4</sup>In contrast,  $I(O; L)$  does not vary significantly, given that the policy does not control which level is played in an episode.

<sup>5</sup>By the data processing inequality we must have  $I(O; \cdot) \geq I(Z; \cdot)$ , and  $C(Z; \cdot)$  cannot be larger than 1. We enforce this upper bound as our estimator sometimes underestimates  $I(O; \cdot)$  for high dimensional observations.

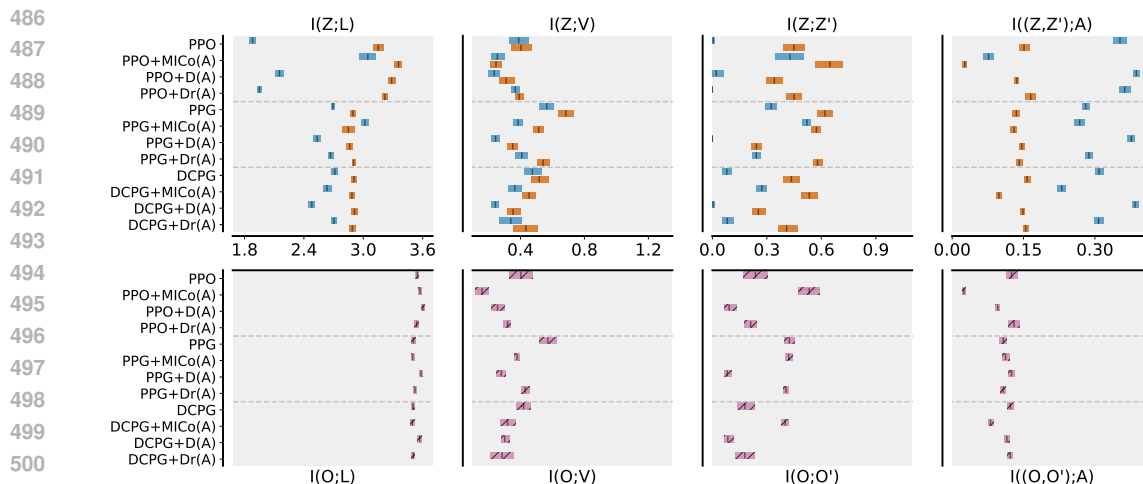


Figure 5: Mean and 95% confidence intervals of  $I(Z; \cdot)/I(O; \cdot)$  (top/bottom) for actor (blue) and critic (orange) representations in Procgen. Information measured from agent observations shown in pink. X-axes are shared across top and bottom. Auxiliary tasks shown are MICO, dynamics prediction (D), and data augmentation (Dr) applied to the actor (A).

## 7 RELATED WORK

**Representation learning in RL.** Representation learning objectives have been used in RL for a variety of reasons such as sample efficiency (Jaderberg et al., 2017; Gelada et al., 2019; Laskin et al., 2020a; Lee et al., 2020; Laskin et al., 2020b), planning (Sekar et al., 2021; McInroe et al., 2024), disentanglement (Dunion et al., 2023), and generalisation (Higgins et al., 2017; Li et al., 2021). Some works focus on designing metrics motivated by theoretical properties such as bisimulation metrics, pseudometrics, decompositions of MDP components, or successor features (Ferns et al., 2004; Mahadevan & Maggioni, 2007; Dayan, 1993; Castro, 2020; Agarwal et al., 2021; Castro et al., 2021; 2023).

**Analysing representations in RL.** Despite the large body of research into representation learning objectives in RL, relatively little work has gone into understanding the learned representations themselves (Wang et al., 2024). Several works use linear probing to determine how well learned representations relate to environment or agent properties (Racah & Pal, 2019; Guo et al., 2019; Anand et al., 2019; Zhang et al., 2024). Other works analyse the learned representation functions via saliency maps which help visualise where an agent is “paying attention” (Rosynski et al., 2020; Atray et al., 2020; Dunion et al., 2024).

## 8 CONCLUSION

In this work, we conducted an in-depth analysis of the representations learned by actor and critic networks in on-policy deep reinforcement learning. Our key findings revealed that when decoupled, actor and critic representations specialise in extracting different types of information from the environment. We found that employing representation learning objectives that support the actor and critic specialisations can result in significant performance gains. Finally, we discovered that the critic influences policy updates to collect data that is informative for its own learning objective. This finding highlighted the critic’s significant role in shaping exploration.

Our work opens up new avenues for research into the interplay between actor and critic representations in reinforcement learning. Future work could explore the implications of our findings for exploration strategies, and whether we observe similar specialisation and interplay outside of the online and on-policy setting.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## REPRODUCIBILITY STATEMENT

Reproducibility can be challenging without access to the data generated during experiments. To assist with this, we will make all of our experimental data, including model checkpoints, logged data and the code for reproducing the figures in this paper openly available upon its publication.

## REFERENCES

- Rishabh Agarwal, Marlos C. Machado, Pablo Samuel Castro, and Marc G. Bellemare. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. In *ICLR*, 2021.
- Cameron Allen, Neev Parikh, Omer Gottesman, and George Konidaris. Learning markov state abstractions for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:8229–8241, 2021.
- Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon Hjelm. Unsupervised state representation learning in atari. 2019.
- Akanksha Atrey, Kaleigh Clary, and David Jensen. Exploratory not explanatory: Counterfactual analysis of saliency maps for deep reinforcement learning. In *ICLR*, 2020.
- Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents (extended abstract). In Qiang Yang and Michael J. Wooldridge (eds.), *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pp. 4148–4152. AAAI Press, 2015. URL <http://ijcai.org/Abstract/15/585>.
- Martín Bertrán, Natalia Martínez, Mariano Phielipp, and Guillermo Sapiro. Instance-based generalization in reinforcement learning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/82674fc29bc0d9895cee346548c2cb5c-Abstract.html>.
- Pablo Samuel Castro. Scalable methods for computing state similarity in deterministic markov decision processes. In *AAAI*, 2020.
- Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, and Mark Rowland. Mico: Improved representations via sampling-based state similarity for markov decision processes. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 30113–30126, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/fd06b8ea02fe5b1c2496fe1700e9d16c-Abstract.html>.
- Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, and Mark Rowland. A kernel perspective on behavioural metrics for markov decision processes. *TMLR*, 2023.
- Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 2048–2056. PMLR, 2020. URL <http://proceedings.mlr.press/v119/cobbe20a.html>.
- Karl Cobbe, Jacob Hilton, Oleg Klimov, and John Schulman. Phasic policy gradient. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 2020–2027. PMLR, 2021. URL <http://proceedings.mlr.press/v139/cobbe21a.html>.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5:613–624, 1993.

- 594 Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.  
595 Provably efficient rl with rich observations via latent state decoding. In *International Conference*  
596 *on Machine Learning*, pp. 1665–1674. PMLR, 2019.
- 597 Mhairi Dunion, Trevor McInroe, Kevin Sebastian Luck, Josiah P. Hanna, and Stefano V. Albrecht.  
598 Temporal disentanglement of representations for improved generalisation in reinforcement learn-  
599 ing. In *ICLR*, 2023.
- 600 Mhairi Dunion, Trevor McInroe, Kevin Sebastian Luck, Josiah Hanna, and Stefano Albrecht. Con-  
601 ditional mutual information for disentangled representations in reinforcement learning. *Advances*  
602 *in Neural Information Processing Systems*, 36, 2024.
- 603 Norman Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision pro-  
604 cesses. In *Conference on Uncertainty in Artificial Intelligence*, 2004.
- 605 C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem.  
606 Brax - a differentiable physics engine for large scale rigid body simulation, 2021. URL <http://github.com/google/brax>.
- 607 Samuel Garcin, James Doran, Shangmin Guo, Christopher G. Lucas, and Stefano V. Albrecht. Dred:  
608 Zero-shot transfer in reinforcement learning via data-regularised environment design. In *Internat-*  
609 *ional Conference on Machine Learning*, 2024. URL <https://api.semanticscholar.org/CorpusID:269449176>.
- 610 Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G. Bellemare. Deepmdp:  
611 Learning continuous latent space models for representation learning. In *ICML*, 2019.
- 612 Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Bernardo A. Pires, and Rémi Munos.  
613 Neural predictive belief representations. *arXiv preprint: arXiv:1811.06407*, 2019.
- 614 Irina Higgins, Arka Pal, Andrei A. Rusu, Loic Matthey, Christopher P Burgess, Alexander Pritzel,  
615 Matthew Botvinick, Charles Blundell, and Alexander Lerchner. Darla: Improving zero-shot trans-  
616 fer in reinforcement learning. In *ICML*, 2017.
- 617 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Ki-  
618 nal Mehta, and João G.M. Araújo. Cleanrl: High-quality single-file implementations of deep  
619 reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.  
620 URL <http://jmlr.org/papers/v23/21-1342.html>.
- 621 Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David  
622 Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv*  
623 *preprint: arXiv:1611.05397*, 2017.
- 624 Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In Marina Meila  
625 and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning,*  
626 *ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning*  
627 *Research*, pp. 4940–4950. PMLR, 2021. URL <http://proceedings.mlr.press/v139/jiang21b.html>.
- 628 Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot gener-  
629 alisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76:201–264,  
630 2023.
- 631 Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Phys-*  
632 *ical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 69(6):066138, 2004.
- 633 Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Re-  
634 inforcement learning with augmented data. In *NeurIPS*, 2020a.
- 635 Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: contrastive unsupervised representa-  
636 tions for reinforcement learning. In *Proceedings of the 37th International Conference on Ma-*  
637 *chine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of*  
638 *Machine Learning Research*, pp. 5639–5650. PMLR, 2020b. URL <http://proceedings.mlr.press/v119/laskin20a.html>.

- 648 Kuang-Huei Lee, Ian Fischer, Anthony Liu, Yijie Guo, Honglak Lee, John Canny, and Sergio  
649 Guadarrama. Predictive information accelerates learning in rl. In *NeurIPS*, 2020.
- 650
- 651 Bonnie Li, Vincent François-Lavet, Thang Doan, and Joelle Pineau. Domain adversarial reinforce-  
652 ment learning. *arXiv preprint: arXiv:2102.07097*, 2021.
- 653 Sridhar Mahadevan and Mauro Maggioni. Proto-value functions: A laplacian framework for learn-  
654 ing representation and control in markov decision processes. *JMLR*, 2007.
- 655
- 656 Trevor McInroe, Lukas Schäfer, and Stefano V. Albrecht. Multi-horizon representations with hier-  
657 archical forward models for reinforcement learning. *TMLR*, 2023.
- 658 Trevor McInroe, Adam Jelley, Stefano V. Albrecht, and Amos Storkey. Planning to go out-of-  
659 distribution in offline-to-online reinforcement learning. In *RLC*, 2024.
- 660
- 661 Seungyong Moon, JunYeong Lee, and Hyun Oh Song. Rethinking value function learning for gen-  
662 eralization in reinforcement learning. *Advances in Neural Information Processing Systems*, 35:  
663 34846–34858, 2022.
- 664 Evan Racah and Christopher Pal. Supervise thyself: Examining self-supervised representations in  
665 interactive environments. *arXiv preprint: arXiv:1906.11951*, 2019.
- 666
- 667 Roberta Raileanu and Rob Fergus. Decoupling value and policy for generalization in reinforce-  
668 ment learning. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th Interna-*  
669 *tional Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, vol-  
670 *ume 139 of Proceedings of Machine Learning Research*, pp. 8787–8798. PMLR, 2021. URL  
671 <http://proceedings.mlr.press/v139/raileanu21a.html>.
- 672 Roberta Raileanu, Max Goldstein, Denis Yarats, Ilya Kostrikov, and Rob Fergus. Auto-  
673 matic data augmentation for generalization in reinforcement learning. In Marc’Aurelio Ran-  
674 zato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan  
675 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-*  
676 *ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.  
677 5402–5415, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/2b38c2df6a49b97f706ec9148ce48d86-Abstract.html)  
678 [2b38c2df6a49b97f706ec9148ce48d86-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/2b38c2df6a49b97f706ec9148ce48d86-Abstract.html).
- 679 Brian C Ross. Mutual information between discrete and continuous data sets. *PLoS one*, 9(2):  
680 e87357, 2014.
- 681 Matthias Rosynski, Frank Kirchner, and Matias Valdenegro-Toro. Are gradient-based saliency maps  
682 useful in deep reinforcement learning? *arXiv preprint: arXiv:2012.01281*, 2020.
- 683
- 684 John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and Pieter Abbeel. High-  
685 dimensional continuous control using generalized advantage estimation. In Yoshua Bengio and  
686 Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016,*  
687 *San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL [http://](http://arxiv.org/abs/1506.02438)  
688 [arxiv.org/abs/1506.02438](http://arxiv.org/abs/1506.02438).
- 689 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
690 optimization algorithms. *arXiv*, 2017.
- 691
- 692 Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bach-  
693 man. Data-efficient reinforcement learning with self-predictive representations. In *ICLR*, 2021.
- 694 Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak.  
695 Planning to explore via self-supervised world models. In *ICML*, 2021.
- 696
- 697 Austin Stone, Oscar Ramirez, Kurt Konolige, and Rico Jonschkowski. The distracting con-  
698 trol suite – a challenging benchmark for reinforcement learning from pixels. *arXiv preprint*  
699 *arXiv:2101.02722*, 2021.
- 700 Han Wang, Erfan Miah, Martha White, Marlos C Machado, Zaheer Abbas, Raksha Kumaraswamy,  
701 Vincent Liu, and Adam White. Investigating the properties of neural network representations in  
reinforcement learning. *Artificial Intelligence*, 330:104100, 2024.

702 Kaixin Wang, Daquan Zhou, Jiashi Feng, and Shie Manno. Ppg reloaded: An empirical study on  
703 what matters in phasic policy gradient. In *ICML*, 2023.

704  
705 Jiayi Weng, Min Lin, Shengyi Huang, Bo Liu, Denys Makoviichuk, Viktor Makoviychuk,  
706 Zichen Liu, Yufan Song, Ting Luo, Yukun Jiang, Zhongwen Xu, and Shuicheng Yan. En-  
707 vPool: A highly parallel reinforcement learning environment execution engine. In S. Koyejo,  
708 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-*  
709 *formation Processing Systems*, volume 35, pp. 22409–22421. Curran Associates, Inc., 2022.  
710 URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/8caaf08e49ddbada6694fae067442ee21-Paper-Datasets_and_Benchmarks.pdf)  
711 [8caaf08e49ddbada6694fae067442ee21-Paper-Datasets\\_and\\_Benchmarks.](https://proceedings.neurips.cc/paper_files/paper/2022/file/8caaf08e49ddbada6694fae067442ee21-Paper-Datasets_and_Benchmarks.pdf)  
712 pdf.

713 Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Im-  
714 proving sample efficiency in model-free reinforcement learning from images. *arXiv preprint:*  
715 *arXiv:1910.01741*, 2019.

716 Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing  
717 deep reinforcement learning from pixels. In *International Conference on Learning Representa-*  
718 *tions*, 2021. URL <https://openreview.net/forum?id=GY6-6sTvGaf>.

719 Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invari-  
720 ant representations for reinforcement learning without reconstruction. In *ICLR*, 2021.

721  
722 Wancong Zhang, Anthony GX-Chen, Vlad Sobal, Yann LeCun, , and Nicolas Carion. Light-weight  
723 probing of unsupervised representations for reinforcement learning. In *RLC*, 2024.

## 724 725 A THEORETICAL RESULTS

726  
727 **Theorem 3.1.** *The difference in returns achieved in train levels and under the full distribution, or*  
728 *generalisation error, has an upper bound that depends on  $I(Z_A; L)$ , with*

$$729 \mathbb{E}_{c \sim \mathcal{U}(L), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] - \mathbb{E}_{c \sim \mathcal{P}(c), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] \leq \sqrt{\frac{2D^2}{|L|}} \times I(Z_A; L), \quad (5)$$

730  
731 where  $c \sim \mathcal{U}(L)$  indicates  $c$  is sampled uniformly over levels in  $L$ ,  $D$  is a constant such that  
732  $|V^\pi(x)| \leq D/2, \forall x, \pi$  and  $Z_A$  is the output space of the actor’s learned representation.

733  
734 *Proof.* This result directly follows from a result obtained by Bertrán et al. (2020) and reproduced  
735 below.

736  
737 **Theorem A.1.** *For any CMDP such that  $|V^\pi(x)| \leq D/2, \forall x, \pi$ , with  $D$  being a constant, then for*  
738 *any set of training levels  $L$ , and policy  $\pi$*

$$739 \mathbb{E}_{c \sim \mathcal{U}(L), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] - \mathbb{E}_{c \sim \mathcal{P}(c), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] \leq \sqrt{\frac{2D^2}{|L|}} \times I(\pi; L), \quad (10)$$

740  
741 Then, as  $\pi = f \circ \phi_A$ , by the data processing inequality we always have  $I(\pi; L) \leq I(Z_A; L)$ , and  
742 therefore,  
743

$$744 \mathbb{E}_{c \sim \mathcal{U}(L), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] - \mathbb{E}_{c \sim \mathcal{P}(c), x_0 \sim \mathcal{P}_0(c)}[V^\pi(x_0)] \leq \sqrt{\frac{2D^2}{|L|}} \times I(\pi; L)$$

$$745 \leq \sqrt{\frac{2D^2}{|L|}} \times I(Z_A; L)$$

746  
747  
748  
749  
750  
751  
752  
753  
754 Garcin et al. (2024) follow the same reasoning and obtain an equivalent result, without restating the  
755 bound. □

**Theorem 3.2.** *if  $\mathcal{T} : \mathbb{X} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{X})$  satisfies the Markov property, and we have  $I((X, X'); A) = I((Z, Z'); A)$  and  $I(X; X') = I(Z; Z')$  for any  $X, X', A, Z, Z'$  collected using policy  $\pi$ , then  $\mathcal{T}_z : \mathbb{Z} \times \mathbb{A} \rightarrow \mathcal{P}(\mathbb{Z})$  satisfies the Markov property when following  $\pi$ .  $\mathcal{T}_z$  always satisfies the Markov property if the above conditions hold for any  $\pi$ .*

*Proof.* This proof has two part. We first demonstrate that the Inverse Model condition of Theorem A.2 from Allen et al. (2021) (reproduced below) is satisfied if and only if  $I((Z, Z'); A) = I((X, X'); A)$ . We then show that if  $I(Z; Z') = I(X; X')$  then the Density Ratio condition is also satisfied.

**Theorem A.2.**  *$\phi$  is a Markov representation if the following conditions hold for every timestep  $t$  and any policy  $\pi$ :*

1. **Inverse Model.** *The inverse dynamic model, defined as  $I(a|s', s) := \frac{\mathcal{T}(s'|a, s)\pi(a|s)}{P^\pi(s'|s)}$ , where  $P^\pi(s'|s) = \sum_{\bar{a} \in \mathbb{A}} \mathcal{T}(s'|\bar{a}, s)\pi(\bar{a}|s)$ , should be equal in the original and reduced MDPs. That is we have  $P^\pi(a|z', z) = P^\pi(a|s, s'), \forall a \in \mathbb{A}, s, s' \in \mathbb{S}$ .*
2. **Density Ratio.** *The original and abstract next-state density ratios are equal when conditioned on the same abstract state:  $\frac{P^\pi(z'|z)}{P^\pi(z')} = \frac{P^\pi(s'|z)}{P^\pi(s')}, \forall x' \in \mathbb{S}$ , where  $P^\pi(s'|z) = \sum_{\bar{s} \in \mathbb{S}} P^\pi(s'|\bar{s})\mu(\bar{s}|z)$  and  $\mu(s|z) = \frac{\mathbf{1}_{\phi(s)=z} P^\pi(s)}{\sum_{\bar{s} \in \mathbb{S}} P^\pi(s|\bar{s})}$ .  $P^\pi(s'|z)$  is the probability of transitioning to state  $s'$  and  $\mu(s|z)$  is the probability of currently being in state  $s$  when in latent state  $z$ .*

We begin with two observations that are useful for our derivation.

Observation A: Given that any  $z \in \mathbb{Z}$  is obtained from the mapping  $x \xrightarrow{\Omega} o \xrightarrow{\phi} z$ , and that  $h = \phi \circ \Omega$  is a deterministic (but not necessarily invertible) function, each element  $x \in \mathbb{X}$  maps to a single element  $z \in \mathbb{Z}$ . It directly follows that  $\forall a, z_1, z_2 \in \mathbb{A} \times \mathbb{Z} \times \mathbb{Z}$ , we have

$$p(a, z_1, z_2) = \sum_{x_1, x_2 \in \mathbb{X}^2} p(a, x_1, x_2) \mathbf{1}[z_1, z_2 = h(x_1), h(x_2)]$$

and

$$p(z_1, z_2) = \sum_{x_1, x_2 \in \mathbb{X}^2} p(x_1, x_2) \mathbf{1}[z_1, z_2 = h(x_1), h(x_2)]$$

Observation B: Let  $P^\pi(a, x, x')$  be the joint distribution of elements in  $(A, X, X')$  collected under policy  $\pi$ , we have  $P^\pi(a, x_1, x_2) > 0$  if and only if  $a, x_1, x_2 \in (A, X, X')$ .

Observation C: Similarly to obs. B, we have  $P^\pi(x_1, x_2) > 0$  if and only if  $x_1, x_2 \in (X, X')$ .

1) *Proving that the Inverse Model condition is satisfied if and only if  $I((Z, Z'); A) = I((X, X'); A)$ .*

The above is equivalent to showing that the Inverse Model condition is satisfied if and only if  $\mathbf{H}(A|Z, Z') = \mathbf{H}(A|X, X')$ . For  $\mathbf{H}(A|Z, Z')$ , we have

$$\mathbf{H}(A|Z, Z') = - \sum_{A, Z, Z'} P^\pi(a, z, z') \log P^\pi(a|z, z')$$

$$\text{(from obs. A)} = - \sum_{\mathbb{A} \times \mathbb{Z} \times \mathbb{Z}} \sum_{x_1, x_2 \in \mathbb{X}^2} P^\pi(a, x_1, x_2) \mathbf{1}[z, z' = h(x_1), h(x_2)] \log P^\pi(a|z, z')$$

$$\text{(from obs. B)} = - \sum_{\mathbb{A} \times \mathbb{X} \times \mathbb{X}} P^\pi(a, x, x') \sum_{\mathbb{Z}^2} \mathbf{1}[z, z' = h(x), h(x')] \log P^\pi(a|z, z')$$

$$= - \sum_{\mathbb{A} \times \mathbb{X} \times \mathbb{X}} P^\pi(a, x, x') \log \prod_{\mathbb{Z}^2} P^\pi(a|z, z')^{\mathbf{1}[z, z' = h(x), h(x)]}$$

$$= - \sum_{\mathbb{A} \times \mathbb{X} \times \mathbb{X}} P^\pi(x, x') P^\pi(a|x, x') \log \prod_{\mathbb{Z}^2} P^\pi(a|z, z')^{\mathbf{1}[z, z' = h(x), h(x)]}$$

$$= - \mathbb{E}_{X, X'} \left[ \sum_{\mathbb{A}} P^\pi(a|x, x') \log \prod_{\mathbb{Z}^2} P^\pi(a|z, z')^{\mathbf{1}[z, z' = h(x), h(x)]} \right]$$

It follows that

$$\begin{aligned} \mathbf{H}(A|Z, Z') - \mathbf{H}(A|X, X') &= \mathbb{E}_{X, X'} \left[ \sum_{\mathbb{A}} P^\pi(a|x, x') \log \frac{P^\pi(a|x, x')}{\prod_{\mathbb{Z}^2} P^\pi(a|z, z') \mathbf{1}_{[z, z'=h(x), h(x')]}]} \right] \\ &= \mathbb{E}_{X, X'} [D_{\text{KL}}(P\|Q)], \end{aligned}$$

with  $P = P^\pi(a|x, x')$  and  $Q = \prod_{z, z' \in \mathbb{Z}, \mathbb{Z}'} P^\pi(a|z, z') \mathbf{1}_{[z, z'=h(x), h(x')]}].$  From Gibbs inequality we always have  $D_{\text{KL}}(p\|q) \geq 0$ , therefore  $\mathbf{I}((Z, Z'); A) = \mathbf{I}((X, X'); A)$  if and only if  $D_{\text{KL}}(P\|Q) = 0 \forall x, x' \in X, X'$ , which is the case if and only if  $P = Q$  almost  $\mu$ -everywhere.

From observation A, any  $x_1, x_2 \in \mathbb{X}^2$  maps to exactly one pair  $z_1, z_2 \in \mathbb{Z}^2$ , and by construction of  $X, X', Z, Z'$ , for any pair  $x, x' \in X, X'$ , we must have  $Q = \prod_{\bar{z}, \bar{z}' \in \mathbb{Z}^2} P^\pi(a|\bar{z}, \bar{z}') \mathbf{1}_{[\bar{z}, \bar{z}'=h(x), h(x')]} = P^\pi(a|z, z')$ , with  $z, z'$  being the corresponding pair in  $Z, Z'$ .

Therefore  $\mathbf{I}((Z, Z'); A) = \mathbf{I}((X, X'); A)$  if and only if  $P^\pi(a|x, x') = P^\pi(a|z, z') \forall x, x', z, z' \in X, X', Z, Z'$ , and we recover the Inverse Model condition.

Conversely, if the Inverse Model condition is not satisfied, then  $\exists x, x', z, z', a \in X, X', Z, Z', A$  for which  $P \neq Q$ . Then  $D_{\text{KL}}(P\|Q) > 0$  at  $x, x'$  and  $\mathbf{I}((Z, Z'); A) < \mathbf{I}((X, X'); A)$ .

2) *Proving that the Density Ratio condition is satisfied if  $\mathbf{I}(Z; Z') = \mathbf{I}(X; X')$ .*

We first show that satisfying

$$\frac{P^\pi(x'|x)}{P^\pi(x')} = \frac{P^\pi(z'|z)}{P^\pi(z')} \quad \forall x, x', z, z' \in X, X', Z, Z' \quad (11)$$

is sufficient for satisfying the Density Ratio condition  $\frac{P^\pi(x'|z)}{P^\pi(x')} = \frac{P^\pi(z'|z)}{P^\pi(z')}$ . We then show that the condition in Equation (11) holds if and only if  $\mathbf{I}(Z; Z') = \mathbf{I}(X; X')$ .

i) *Showing the Density Ratio condition holds when Equation (11) is satisfied.* First we notice that,  $\forall x', z \in X', Z$ , we have

$$P^\pi(x'|z) = \sum_{\bar{x} \in \mathbb{X}} \mathbf{1}[z = h(\bar{x})] P^\pi(x'|\bar{x}) = \mathbb{E}_X [P^\pi(x'|x)].$$

Then, supposing Equation (11) holds, we must have

$$P^\pi(x'|z) = \mathbb{E}_X [P^\pi(x'|x)] = P^\pi(x') \frac{P^\pi(z'|z)}{P^\pi(z')} \quad \forall x', z, z' \in X', Z, Z',$$

and the Density Ratio condition holds.

ii) *Proving Equation (11) holds if and only if  $\mathbf{I}(Z; Z') = \mathbf{I}(X; X')$ .*

We have

$$\begin{aligned} \mathbf{I}(Z; Z') &= \sum_{\mathbb{Z}^2} P^\pi(z, z') \log \frac{P^\pi(z'|z)}{P^\pi(z')} \\ \text{(from obs. A)} &= \sum_{\mathbb{Z}^2} \sum_{x_1, x_2 \in \mathbb{X}^2} P^\pi(x_1, x_2) \mathbf{1}[z, z' = h(x_1), h(x_2)] \log \frac{P^\pi(z'|z)}{P^\pi(z')} \\ \text{(from obs. C)} &= \sum_{\mathbb{X}^2} P^\pi(x, x') \sum_{\mathbb{Z}^2} \mathbf{1}[z, z' = h(x), h(x')] \log \frac{P^\pi(z'|z)}{P^\pi(z')} \\ &= \mathbb{E}_{X, X'} \left[ \log \prod_{\mathbb{Z}^2} \left( \frac{P^\pi(z'|z)}{P^\pi(z')} \right)^{\mathbf{1}_{[z, z'=h(x), h(x')]}]} \right]. \end{aligned}$$

Then,

$$\mathbf{I}(X; X') - \mathbf{I}(Z; Z') = \mathbb{E}_{X, X'} [D_{\text{KL}}(P'\|Q)],$$

with

$$P' = \frac{P^\pi(x'|x)}{P^\pi(x')} \quad \text{and} \quad Q = \prod_{\mathbb{Z}^2} \left( \frac{P^\pi(z'|z)}{P^\pi(z')} \right)^{\mathbf{1}_{[z, z'=h(x), h(x')]}}$$



The remainder of this part follows the same structure as for the first part of the proof.

$I(X; X') = I(Z; Z')$  if and only if  $\forall x, x' \in X, X', P = Q$  almost  $\mu$ -everywhere. Any  $x_1, x_2 \in \mathbb{X}^2$  maps to exactly one pair  $z_1, z_2 \in \mathbb{Z}^2$ , and by construction of  $X, X', Z, Z'$ , for any pair  $x, x' \in X, X'$ , we must have

$$Q = \prod_{\bar{z}, \bar{z}' \in \mathbb{Z}^2} \left( \frac{P^\pi(\bar{z}'|\bar{z})}{P^\pi(\bar{z}')} \right)^{\mathbf{1}[\bar{z}, \bar{z}' = h(x), h(x')]} = \frac{P^\pi(z'|z)}{P^\pi(z')},$$

with  $z, z'$  being the corresponding pair in  $Z, Z'$ .

Therefore  $I(X; X') = I(Z; Z')$  if and only if  $\forall x, x', z, z' \in X, X', Z, Z'$  we have  $\frac{P^\pi(x'|x)}{P^\pi(x')} = \frac{P^\pi(z'|z)}{P^\pi(z')}$ . Finally, from i) being true, the Density ratio condition must hold.  $\square$

**Lemma 4.1.**  $I(Z; L) > 0$  if  $\exists z_k, c_j \in Z \times L$  such that  $\mu(z_k|c_j) \neq \mu(z_k)$  and  $I(O; L) > 0$ ,  $I(O; L) > 0$  being the mutual information between  $L$  and observations  $O$ , with  $\phi(o) = z \in Z$ .

*Proof.* Given  $\pi$  is fixed while the batch  $O$  is collected, for a single batch the causal interaction between  $L, O$  and  $Z$  is described by the Markov chain  $X \rightarrow O \rightarrow Z$ , where  $x = (s, c) \in \mathbb{S} \times L$  and isn't directly observed. By the data processing inequality,  $I(L; Z) \leq I(L; O)$ , and as such  $I(L; O) > 0$  is a necessary condition for  $I(L; Z)$  to be positive.

Note that

$$I(L; Z) = \mathbf{H}(L) + \mathbf{H}(Z) - \mathbf{H}(L, Z) = 0 \Leftrightarrow \mathbf{H}(L, Z) = \mathbf{H}(L) + \mathbf{H}(Z),$$

that is, if and only if  $Z$  and  $L$  are independently distributed. Given the causal relationship between  $L$  and  $Z$ ,  $\mu(z|c)$  is well defined  $\forall z, c \in Z \times L$ . If  $\exists z_k, c_j \in Z \times L$  such that  $\mu(z_k|c_j) \neq \mu(z_k)$  then  $Z$  and  $L$  cannot be independently distributed, and  $I(L; Z) > 0$ .  $\square$

**Lemma 4.2.**  $I(Z; L)$  monotonically increases with  $I(Z; Z') - I(Z; Z'|L)$ .

*Proof.* Consider an episode of arbitrary length  $N$  collected with policy  $\pi$ . We depict the causal structure that exists between elements in the top row of Figure 6 (elements may be repeated within each sequence). It naturally follows that we have the causal structure depicted in the bottom row when considering all levels in  $L$ . By the chain rule of mutual information, we have

$$I(Z; L) = I(Z; (Z', L)) - I(Z; Z'|L) = I(Z; L|Z') + I(Z; Z') - I(Z; Z'|L),$$

and it follows that  $I(Z; L)$  increases with  $I(Z; Z') - I(Z; Z'|L)$ .

Note that  $I(Z; Z'|L)$  quantifies the dependency between  $Z$  and  $Z'$  that exists regardless of their shared context  $c$ .  $I(Z; Z'|L) = I(Z; Z')$  implies that  $(Z, Z')$  and  $L$  are independent and latent transitions are invariant to the training level.  $\square$

**Lemma 4.3.**  $I(Z; V) > 0$  if  $\exists z_k, v_m \in Z \times V$  such that  $\frac{1}{L} \sum_{c \in L} p(z_k, v_m|c) \neq p(z_k)p(v_m)$ .

*Proof.* We have  $I(Z; V) = 0$  if and only if  $Z$  and  $V$  are independently distributed. If  $\frac{1}{L} \sum_{c \in L} p(z_k, v_m|c) \neq p(z_k)p(v_m)$  for some  $z_k, v_m \in Z \times V$ , then  $p(z, v) = \frac{1}{L} \sum_{c \in L} p(z, v|c) \neq p(z)p(v)$  and  $L$  and  $V$  cannot be independent.  $\square$

**Corollary 4.4.**  $I(Z; V)$  can be positive when  $Z|L$  and  $V|L$  are conditionally independent. If  $I(Z; V) > 0$  and  $Z|L$  and  $V|L$  are conditionally independent, then  $I(Z; L) > 0$ .

*Proof.* Since  $V$  depends on the states and CMDP dynamics given a fixed  $\pi$ , and when  $Z|L$  and  $V|L$  are conditionally independent, we have the causal structure  $V \leftarrow X \rightarrow O \rightarrow Z$ , with no direct causal link between  $Z$  and  $V$ .

However conditional independence does not guarantee independence. If  $\sum_{c \in L} p(z_k|c)p(v_m|c) \neq \sum_{c \in L} p(z_k|c) \sum_{c \in L} p(v_m|c)$  for some  $z_k, v_m \in Z \times V$ , then we still would have  $I(Z; V) > 0$ .

Given the causal structure, and by the data processing inequality,  $I(Z; V) > 0$  directly implies that  $I(Z; L) > 0$ .  $\square$

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

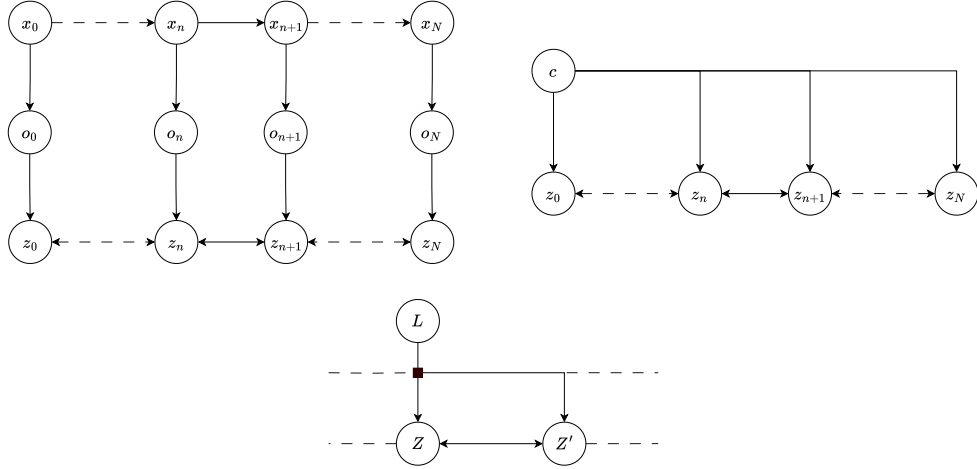


Figure 6: In the top row, left, we depict the causal graph of states, observation and latents obtained over an episode. On the same row we draw a simplified graph that focuses on the relationship between  $c$  and  $Z_{0:N}$ , and utilises the notion that the context remains the same throughout the episode. In the bottom row we draw the resulting causal relationship between  $L$ ,  $Z$  and  $Z'$ .

**Corollary 4.5.**  $I((Z; Z'); A)$  can still be positive when  $(Z, Z')|L$  and  $A|L$  are conditionally independent. If  $I((Z; Z'); A) > 0$  and  $(Z, Z')|L, A|L$  are conditionally independent, then  $\{I(Z; L) > 0$  and  $I(A; L) > 0\}$  and/or  $\{I(Z'; L) > 0$  and  $I(A; L) > 0\}$ .

*Proof.*  $I((Z, Z'); A) \geq \max(I(Z; A), I(Z'; A))$ , therefore showing that either is  $I(Z; A)$  or  $I(Z'; A)$  is positive is sufficient.

The rest of the proof directly follows the proof for Corollary 4.4. We will have  $I(Z; A) > 0$ ,  $I(Z; L) > 0$  and  $I(A; L) > 0$  given the causal chain  $A \leftarrow X \rightarrow O \rightarrow Z$ , even under conditional independence between  $Z|L$  and  $A|L$ . Similarly, we will have  $I(Z'; A) > 0$ ,  $I(Z'; L) > 0$  and  $I(A; L) > 0$  given the causal chain  $A \leftarrow X' \rightarrow O' \rightarrow Z'$ .  $\square$

## B ADDITIONAL FIGURES AND TABLES

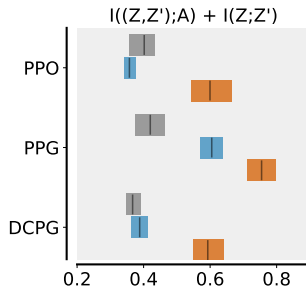


Figure 7:  $I((Z, Z'); A) + I(Z; Z')$  for shared (gray), actor (blue) and critic (orange) for PPO, PPG, and DCPG in Procgen.

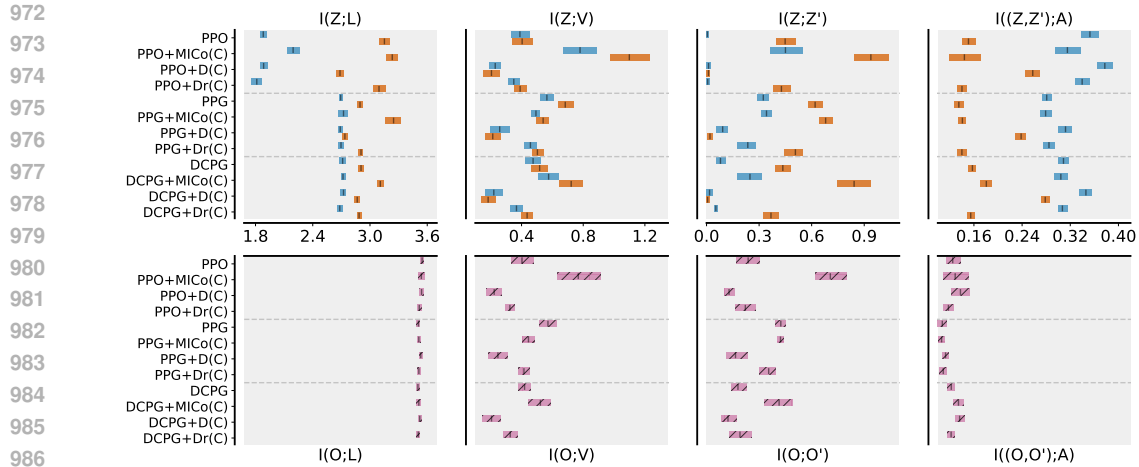


Figure 8: Mutual information measurements for the actor (blue) and critic (orange) for auxiliary losses applied to the critic for PPO, PPG, and DCPG in Procgen. Top/bottom rows are  $I(Z; \cdot)/I(O; \cdot)$  with a shared x-axis.

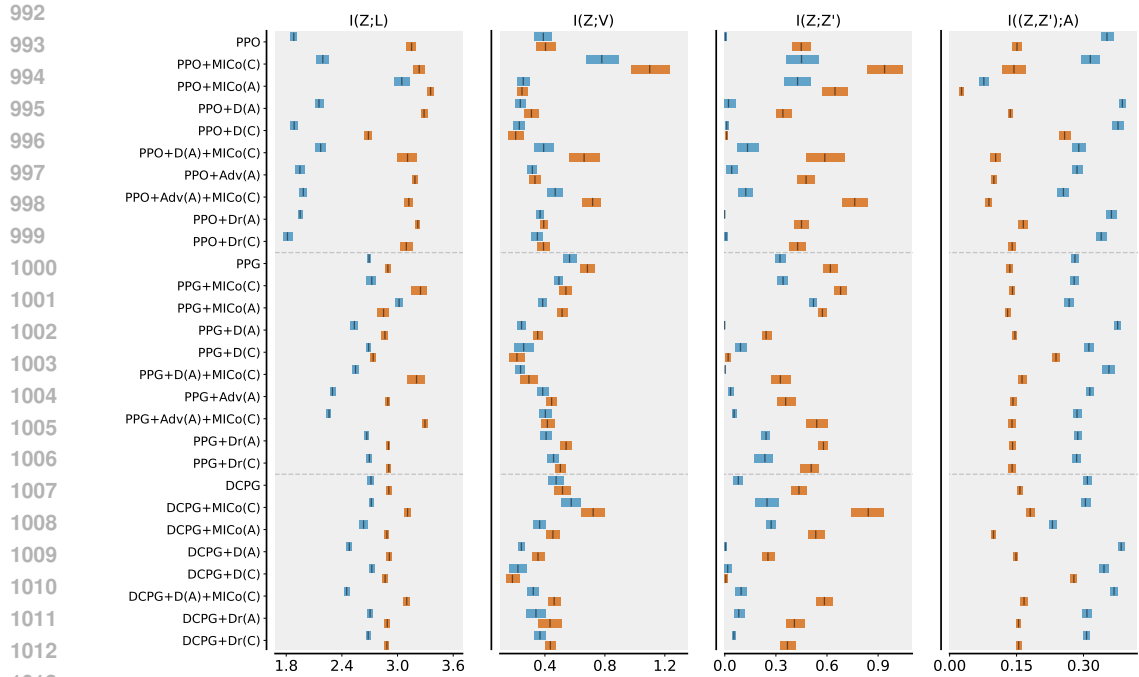


Figure 9:  $I(Z; \cdot)$  measurements for the actor (blue) and critic (orange) for auxiliary losses for PPO, PPG, and DCPG in Procgen.

## C IMPLEMENTATION DETAILS

### C.1 MUTUAL INFORMATION ESTIMATION

We measure mutual information using the estimator proposed by Kraskov et al. (2004) and later extended to pairings of continuous and discrete variables by Ross (2014). These methods are based on performing entropy estimation using  $k$ -nearest neighbors distances. We use  $k = 3$  and determine nearest neighbors by measuring the Euclidian ( $L_2$ ) distance between points. We checked measure-

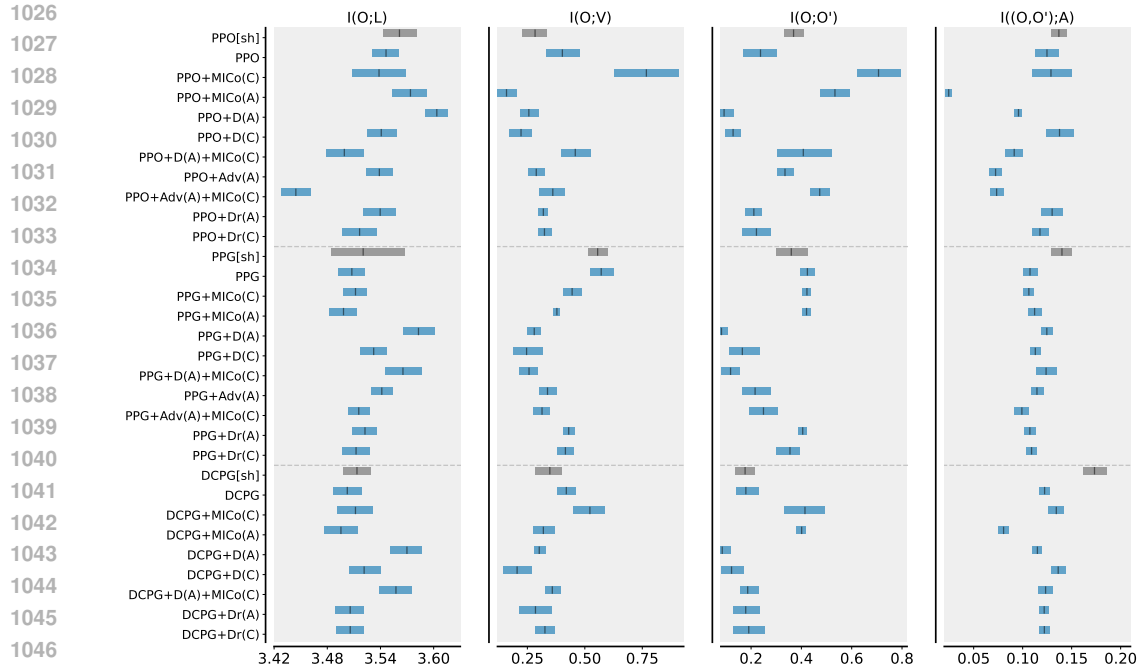


Figure 10:  $I(O; \cdot)$  measurements for the actor (blue) and critic (orange) for auxiliary losses for PPO, PPG, and DCPG in Procgen.

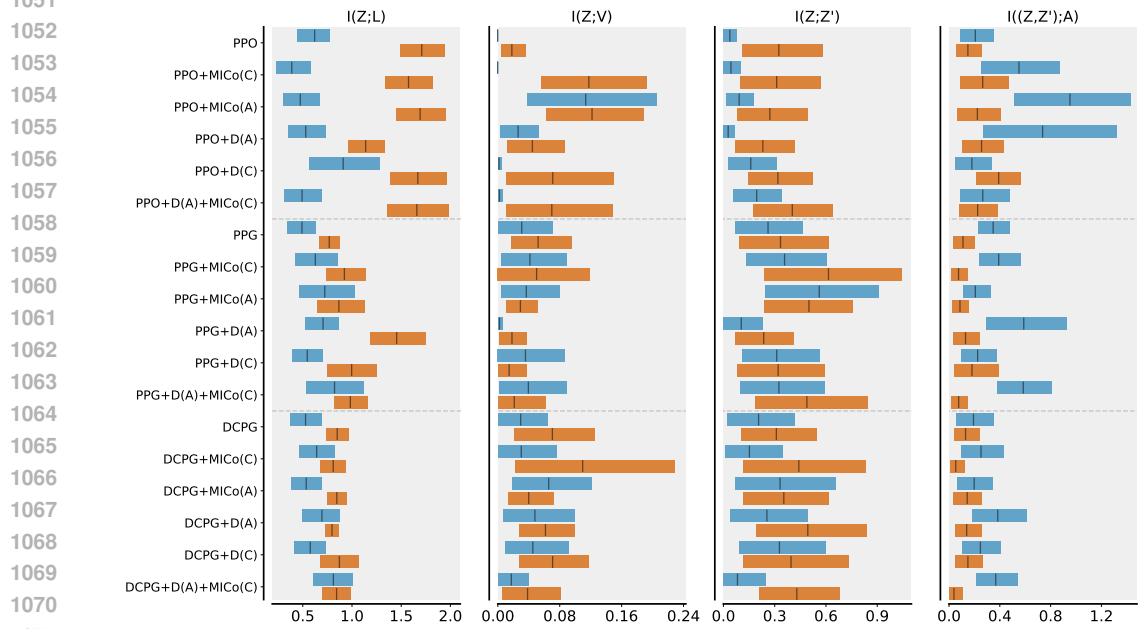


Figure 11: Mutual information measurements for the actor (blue) and critic (orange) for auxiliary losses for PPO, PPG, and DCPG in Brax.

ments obtained when using different  $k$  and under different metric spaces, and we found that our measurements are broadly invariant to the choice of estimator parameters.

At the end of training we collect a batch of trajectories consisting of  $2^{16}$  timesteps ( $2^{15}$  timesteps in Brax) from  $L$ . We construct  $(A, O, O', Z, Z', V, L)$  from  $n = 4096$  timesteps yielding

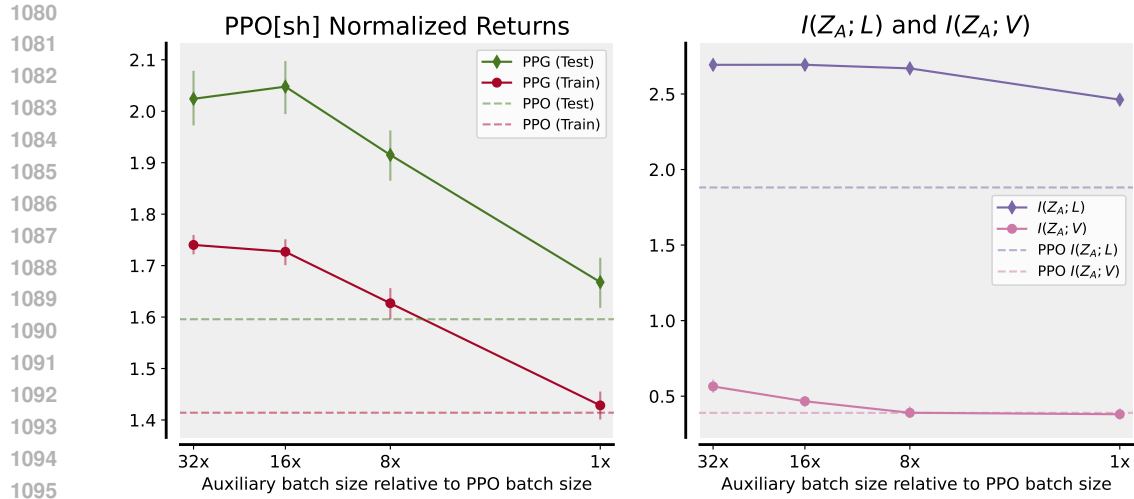


Figure 12: Progen PPG returns (left) normalized by PPO[sh] performance and mutual information quantities  $I(Z_A; L)/I(Z_A; V)$  (right) for varying auxiliary batch size levels.

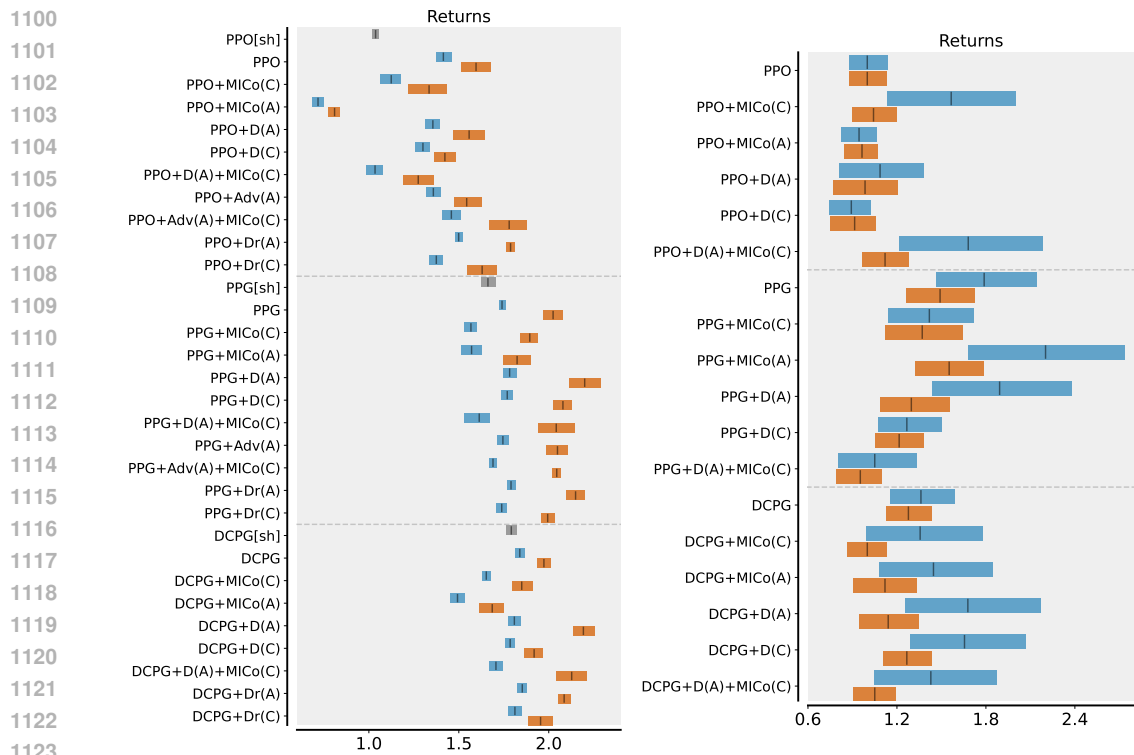


Figure 13: Returns in Progen (left) and Brax (right).

$(a_t, o_t, o_{t+1}, z_t, z_{t+1}, v_t, c_t)$ . Subsampling is necessary to compute mutual information estimates in a reasonable time, while ensuring we sample states from most levels in  $L$  and at various point of the trajectories followed by the agent in each level. Timesteps are sampled uniformly and without replacement from the batch, after having excluded:

1. Odd timesteps, to ensure  $O$  and  $O'$  will not overlap (i.e.  $O$  contains only even timesteps, and  $O'$ , being sampled at  $t + 1$ , contains only odd timesteps).

1134 Table 2: Measurements of compression efficiency  $C(Z_A|O; V)$  (Equation (9)) with standard error  
 1135 in Procgen. Statistical significance bolded, determined by Welch’s t-test. Results highlighted in  
 1136 red when decoupling decreases  $C(Z_A|O; V)$ , and highlighted in green when decoupling increases  
 1137  $C(Z_A|O; V)$ , otherwise yellow. Coupled architectures are denoted with algorithm name plus “[sh]”.  
 1138

Algorithm	$C(Z_A O; V)$	$C(Z_A O; L)$
PPO[sh]	89.3 ± 2	65.2 ± 3
PPO	90.1 ± 4	52.3 ± 3
PPG[sh]	85.9 ± 4	70.0 ± 2
PPG	94.1 ± 2	75.5 ± 2
DCPG[sh]	95.4 ± 2	77.6 ± 2
DCPG	92.3 ± 7	76.4 ± 2

1147 Table 3: Measurements of compression efficiency  $C(Z_A|O; \cdot)$  (Equation (9)) of the actor’s repre-  
 1148 sentation  $\phi_A$  in Procgen. Results highlighted in red when the auxiliary loss decreases the metric  
 1149 relative to the base algorithm, and highlighted in green when the auxiliary loss increases the metric  
 1150 relative to the base algorithm. Auxiliary losses are applied to the actor (A) and critic (C) in the form  
 1151 of dynamics prediction (D), MICo, and advantage distillation (Adv).  
 1152

Algorithm	$C(Z_A O; V)$	$C(Z_A O; L)$	$C((Z_A O, Z'_A O'); A)$
PPO	90.1 ± 4	52.3 ± 3	99.9 ± 0
PPO+MICo(C)	93.9 ± 2	60.4 ± 3	99.4 ± 0
PPO+MICo(A)	98.6 ± 1	84.7 ± 2	87.5 ± 5
PPO+D(A)	62.6 ± 12	61.3 ± 2	100.0 ± 0
PPO+D(C)	86.8 ± 7	52.8 ± 3	100.0 ± 0
PPO+D(A)+MICo(C)	76.3 ± 6	63.7 ± 3	99.5 ± 0
PPO+Adv(A)	96.9 ± 2	53.2 ± 3	100.0 ± 0
PPO+Adv(A)+MICo(C)	100.0 ± 0	57.0 ± 2	98.5 ± 1
PPO+Dr(A)	89.1 ± 6	54.3 ± 3	100.0 ± 0
PPO+Dr(C)	98.1 ± 1	50.8 ± 3	99.6 ± 0
PPG	94.1 ± 2	75.5 ± 2	100.0 ± 0
PPG+MICo(C)	98.0 ± 1	76.4 ± 2	100.0 ± 0
PPG+MICo(A)	95.4 ± 2	85.8 ± 2	100.0 ± 0
PPG+D(A)	89.5 ± 4	71.2 ± 2	100.0 ± 0
PPG+D(C)	96.3 ± 2	75.1 ± 2	100.0 ± 0
PPG+D(A)+MICo(C)	85.9 ± 7	71.6 ± 2	100.0 ± 0
PPG+Adv(A)	98.4 ± 1	63.3 ± 2	100.0 ± 0
PPG+Adv(A)+MICo(C)	99.4 ± 1	62.9 ± 2	100.0 ± 0
PPG+Dr(A)	91.3 ± 3	74.9 ± 2	100.0 ± 0
PPG+Dr(C)	93.1 ± 6	75.6 ± 2	100.0 ± 0
DCPG	92.3 ± 7	76.4 ± 2	100.0 ± 0
DCPG+MICo(C)	98.1 ± 1	76.6 ± 2	100.0 ± 0
DCPG+MICo(A)	91.7 ± 3	74.3 ± 2	100.0 ± 0
DCPG+D(A)	80.9 ± 4	69.9 ± 2	100.0 ± 0
DCPG+D(C)	97.8 ± 1	76.3 ± 2	100.0 ± 0
DCPG+D(A)+MICo(C)	83.4 ± 5	69.6 ± 2	100.0 ± 0
DCPG+Dr(A)	97.5 ± 2	76.7 ± 2	100.0 ± 0

1185  
 1186  
 1187 2. Timesteps corresponding to episode terminations, to ensure the pair  $o_t, o_{t+1}$  cannot orig-  
 inate from different levels.

Table 4: Measurements of compression efficiency  $C(Z_C|O; \cdot)$  (Equation (9)) of the actor’s representation  $\phi_C$  in Procgen. Results highlighted in red when the auxiliary loss decreases the metric relative to the base algorithm, highlighted in green when the auxiliary loss increases the metric relative to the base algorithm, and highlighted in yellow otherwise. Auxiliary losses are applied to the actor (A) and critic (C) in the form of dynamics prediction (D), MICo, and advantage distillation (Adv).

Algorithm	$C(Z_C O; V)$	$C(Z_C O; L)$	$C((Z_C O, Z'_C O'; A))$
PPO	93.7 ± 3	88.4 ± 2	85.6 ± 4
PPO+MICo(C)	100.0 ± 0	90.3 ± 1	82.7 ± 3
PPO+MICo(A)	97.6 ± 2	92.4 ± 1	64.4 ± 6
PPO+D(A)	99.7 ± 0	90.2 ± 1	87.4 ± 3
PPO+D(C)	87.6 ± 4	77.4 ± 2	99.1 ± 0
PPO+D(A)+MICo(C)	91.0 ± 6	88.1 ± 2	84.4 ± 3
PPO+Adv(A)	96.7 ± 2	89.3 ± 1	87.6 ± 4
PPO+Adv(A)+MICo(C)	100.0 ± 0	89.9 ± 1	87.0 ± 3
PPO+Dr(A)	98.0 ± 1	90.0 ± 1	87.9 ± 3
PPO+Dr(C)	97.5 ± 1	87.0 ± 2	86.3 ± 3
PPG	99.2 ± 1	81.6 ± 2	90.3 ± 3
PPG+MICo(C)	93.0 ± 6	90.9 ± 2	91.8 ± 2
PPG+MICo(A)	100.0 ± 0	80.9 ± 2	84.4 ± 4
PPG+D(A)	100.0 ± 0	79.2 ± 2	91.3 ± 3
PPG+D(C)	89.4 ± 4	77.6 ± 2	100.0 ± 0
PPG+D(A)+MICo(C)	93.3 ± 4	89.1 ± 2	87.1 ± 4
PPG+Adv(A)	100.0 ± 0	80.9 ± 2	89.7 ± 3
PPG+Adv(A)+MICo(C)	99.9 ± 0	92.3 ± 1	93.1 ± 3
PPG+Dr(A)	98.7 ± 1	81.7 ± 2	90.2 ± 3
PPG+Dr(C)	96.4 ± 3	81.8 ± 2	85.8 ± 4
DCPG	99.5 ± 1	81.7 ± 2	92.1 ± 3
DCPG+MICo(C)	98.9 ± 1	88.6 ± 2	93.5 ± 2
DCPG+MICo(A)	99.8 ± 0	81.4 ± 2	87.2 ± 4
DCPG+D(A)	100.0 ± 0	80.7 ± 2	92.6 ± 3
DCPG+D(C)	91.3 ± 3	81.2 ± 1	100.0 ± 0
DCPG+D(A)+MICo(C)	99.6 ± 0	87.4 ± 2	92.7 ± 3
DCPG+Dr(A)	100.0 ± 0	81.3 ± 2	89.4 ± 3
DCPG+Dr(C)	97.0 ± 2	81.2 ± 2	90.7 ± 3

3. Timesteps from episodes that have not terminated, to ensure we can always compute  $v_t$ .

## C.2 PROCGEN

The Procgen Benchmark is a set of 16 diverse PCG environments that echoes the gameplay variety seen in the ALE benchmark Bellemare et al. (2015). The game levels, determined by a random seed, can differ in visual design, navigational structure, and the starting locations of entities. All Procgen environments use a common discrete 15-dimensional action space and generate  $64 \times 64 \times 3$  RGB observations. A detailed description of each of the 16 environments is provided by Cobbe et al. (2020). RL algorithms such as PPO reveal significant differences between test and training performance in all games, making Procgen a valuable tool for evaluating generalisation performance.

We conduct our experiment on the easy setting of Procgen, which employs 200 training levels and a budget of 25M training steps, and evaluate the agent’s scores on the training levels and on the full range of levels, excluding the training levels. We use the version of Procgen provided by EnvPool

1242 (Weng et al., 2022). Following prior work, (Raileanu et al., 2021; Jiang et al., 2021; Moon et al.,  
1243 2022), for each game we normalise train/test scores by the mean train/test score achieved by PPO  
1244 in that game.

1245 For PPO, we base our implementation on the CleanRL PPO implementation (Huang et al., 2022),  
1246 which reimplements the PPO agent from the original Procgen publication in JAX. We use the same  
1247 ResNet policy architecture and PPO hyperparameters (identical for all games) as Cobbe et al. (2020)  
1248 and reported in Table 5.

1249 We re-implement PPG and DCPG in JAX, based on the Pytorch implementations provided by Huang  
1250 et al. (2022) and Moon et al. (2022). We use the default recommended hyperparameters for each  
1251 algorithms, which are reported in Table 6. We note that our PPG implementation ends up outper-  
1252 forming the original implementation by about 10% on the test set, while our DCPG implementation  
1253 underperforms test scores reported by Moon et al. (2022) by about 10%. We attribute this discrep-  
1254 ancy to minor differences between the JAX and Pytorch libraries, and decided to not investigate  
1255 further.

1256 We conduct our experiments on A100 and RTX8000 Nvidia GPUs and 6 CPU cores. One seed for  
1257 one game completes in 2 to 12 hours, depending on the GPU, algorithm, and whether the architecture  
1258 is coupled or decoupled (for example, PPG decoupled can be expected to run 4x to 6x slower than  
1259 PPO coupled).

### 1261 C.3 BRAX

1262  
1263 For our experiments in Brax, we implement a custom “video distractors” set of tasks, similar to  
1264 those from (Stone et al., 2021). In this setup, a video plays in an overlay on the pixels the agent  
1265 views. There is a disjoint set of videos between the training and testing environments. The random  
1266 seed determines the environment’s initial physics and the video overlay at the beginning of training.  
1267 The pixels themselves are full-RGB  $64 \times 64 \times 3$  arrays, but we use framestacking to bring each  
1268 agent input to  $64 \times 64 \times 9$  pixels.

1269 Similar to the algorithms used in the Procgen experiments, we implement our algorithms in JAX  
1270 and base them on CleanRL.

1271 We conduct our experiments on RTX A4500 Nvidia GPUs and 6 CPU cores. One seed completes  
1272 in 7.5-48 hours, depending on the environment and its physics backend as well as the algorithm.

1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295



1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

Table 5: Hyperparameters used for PPO in Procgen and Brax experiments. All runs employing a specific (or combination of) representation learning objective use the same hyperparameters.

Parameter	Procgen	Brax
<i>PPO</i>		
$\gamma$	0.999	0.999
$\lambda_{\text{GAE}}$	0.95	0.95
rollout length	256	128
minibatches per epoch	8	8
minibatch size	2048	512
$J_\pi$ clip range	0.2	0.2
number of environments	64	32
Adam learning rate	5e-4	5e-4
Adam $\epsilon$	1e-5	1e-8
max gradient norm	0.5	0.5
value clipping	no	no
return normalisation	yes	no
value loss coefficient	0.5	0.5
entropy coefficient	0.01	0.01
<i>PPO (coupled)</i>		
PPO epochs (actor and critic)	3	-
<i>PPO (decoupled)</i>		
Actor epochs	1	1
Critic epochs	9	1
<i>MICo objective</i>		
MICo coefficient	0.5	0.01
Target network update coefficient	0.005	0.05
<i>Dynamics objective</i>		
Dynamics loss coefficient	1.0	0.01
In-distribution transitions weighting	1.0	1.0
Out-of-distribution states weighting	1.0	1.0
Out-of-distribution actions weighting	0.5	0.5
<i>Advantage distillation objective</i>		
Advantage prediction coefficient	0.25	-

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

Table 6: Hyperparameters used for PPG and DCPG in Procgen experiments. Hyperparameters shared between methods are only reported if they change from the method above. All runs employing a specific (or combination of) representation learning objective use the same hyperparameters.

Parameter	Procgen
<i>PPG</i>	
$\gamma$	0.999
$\lambda_{\text{GAE}}$	0.95
rollout length	256
minibatches per epoch policy phase	8
minibatch size policy phase	2048
minibatches per epoch auxiliary phase	512
minibatch size auxiliary phase	1024
$J_\pi$ clip range	0.2
number of environments	64
Adam learning rate	5e-4
Adam $\epsilon$	1e-5
max gradient norm	0.5
value clipping	no
return normalisation	yes
value loss coefficient policy phase	0.5
value loss coefficient auxiliary phase	1.0
entropy coefficient	0.01
policy phase epochs	1
auxiliary phase epochs	6
number of policy phases per auxiliary phase	32
policy regularisation coefficient $\beta_c$	1.0
auxiliary value distillation coefficient	1.0
<i>DCPG</i>	
value loss coefficient policy phase	0.0
delayed value loss coefficient policy phase	1.0
<i>MICo objective</i>	
MICo coefficient	0.5
Target network update coefficient	0.005
<i>Dynamics objective</i>	
Dynamics loss coefficient	1.0
In-distribution transitions weighting	1.0
Out-of-distribution states weighting	1.0
Out-of-distribution actions weighting	0.5