

---

# Decoupling Value and Policy for Generalization in Reinforcement Learning

---

Roberta Raileanu<sup>1</sup> Rob Fergus<sup>1</sup>

## Abstract

Standard deep reinforcement learning algorithms use a shared representation for the policy and value function, especially when training directly from images. However, we argue that more information is needed to accurately estimate the value function than to learn the optimal policy. Consequently, the use of a shared representation for the policy and value function can lead to overfitting. To alleviate this problem, we propose two approaches which are combined to create IDAAC: Invariant Decoupled Advantage Actor-Critic. First, IDAAC decouples the optimization of the policy and value function, using separate networks to model them. Second, it introduces an auxiliary loss which encourages the representation to be invariant to task-irrelevant properties of the environment. IDAAC shows good generalization to unseen environments, achieving a new state-of-the-art on the Procgen benchmark and outperforming popular methods on DeepMind Control tasks with distractors. Our implementation is available at <https://github.com/rraileanu/idaac>.

## 1. Introduction

Generalization remains one of the main challenges of deep reinforcement learning (RL). Current methods fail to generalize to new scenarios even when trained on semantically similar environments with the same high-level goal but different dynamics, layouts, and visual appearances (Farebrother et al., 2018; Packer et al., 2018; Zhang et al., 2018a; Cobbe et al., 2018; Gamrian & Goldberg, 2019; Cobbe et al., 2019; Song et al., 2020). This indicates that standard RL agents memorize specific trajectories rather than learning transferable skills. Several strategies have been proposed to alleviate this problem, such as the use of regular-

ization (Farebrother et al., 2018; Zhang et al., 2018a; Cobbe et al., 2018; Igl et al., 2019), data augmentation (Cobbe et al., 2018; Lee et al., 2020; Ye et al., 2020; Kostrikov et al., 2020; Laskin et al., 2020; Raileanu et al., 2020), or representation learning (Zhang et al., 2020a;c; Mazouze et al., 2020; Stooke et al., 2020; Agarwal et al., 2021).

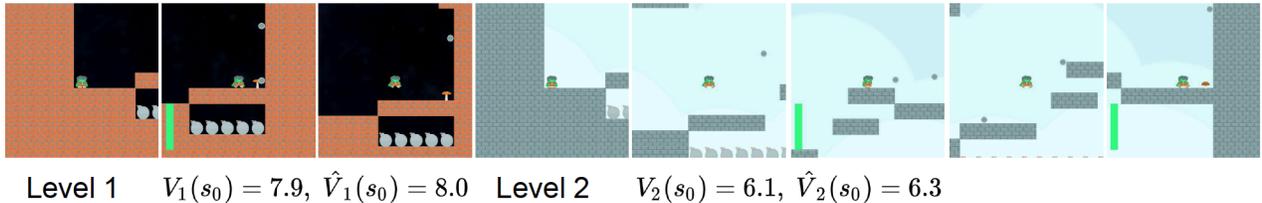
Here we consider the problem of generalizing to unseen instances (or levels) of procedurally generated environments, after training on a relatively small number of such instances. While the high-level goal is the same, the background, dynamics, layouts, as well as the locations, shapes, and colors of various entities, differ across instances.

In this work, we identify a new factor that leads to overfitting in such settings, namely the use of a shared representation for the policy and value function. We point out that accurately estimating the value function requires instance-specific features in addition to the information needed to learn the optimal policy. When training a common network for the policy and value function (as is currently standard practice in pixel-based RL), the need for capturing level-specific features in order to estimate the value function can result in a policy that does not generalize well to new task instances.

### 1.1. Policy-Value Representation Asymmetry

To illustrate this phenomenon, which we call the *policy-value representation asymmetry*, consider the example in Figure 1 which shows two different levels from the Procgen game Ninja (Cobbe et al., 2019). The first observations of the two levels are semantically identical and could be represented using the same features (describing the locations of the agent, the bombs, and the platform, while disregarding the background patterns and wall colors). An optimal agent should take the same action in both levels, namely that of moving to the right to further explore the level and move closer to the goal (which is always to its right in this game). However, these two observations have different values. Note that Level 1 is much shorter than Level 2, and the levels look quite different from each other after the initial part. A standard RL agent (*e.g.* PPO (Schulman et al., 2017)) completes the levels in 24 and 50 steps, respectively. The true value of the initial observation is the expected return (*i.e.* sum of discounted rewards) received during the episode. In

<sup>1</sup>Department of Computer Science, New York University, New York, USA. Correspondence to: Roberta Raileanu <rraileanu@cs.nyu.edu>.



**Figure 1. Policy-Value Asymmetry.** Two Ninja levels with initial observations that are *semantically identical but visually different*. Level 1 (first three frames from the left with black background) is much shorter than Level 2 (last five frames with blue background). Both the true values and the estimated values (by a PPO agent trained on 200 levels) of the initial observation are higher for Level 1 than for Level 2 *i.e.*  $V_1(s_0) > V_2(s_0)$  and  $\hat{V}_1(s_0) > \hat{V}_2(s_0)$ . Thus to accurately predict the value function, the representations must capture level-specific features (such as the backgrounds), which are irrelevant for learning the optimal policy. Consequently, using a common representation for both the policy and value function can lead to overfitting to spurious correlations and poor generalization to unseen levels.

this game, the agent receives a reward of 10 when it reaches the goal and 0 otherwise. Hence, the true value is higher for Level 1 than for Level 2 since the reward is discounted only for 24 steps rather than 50. In order to accurately estimate the value of an observation (which is part of the objective function for many popular RL methods), the agent must memorize the number of remaining steps in that particular level (which for the initial observation is equivalent to the episode’s length). To do this, the agent must use instance-specific features such as the background (which can vary across a level so that each observation within a level has a slightly different background pattern). For the case shown in Figure 1, the agent can learn to associate a black background with a higher value than a blue background. But this is only a spurious correlation since the background has no causal effect on the state’s value, unlike the agent’s position relative to items it can interact with. If an agent uses a shared representation for learning the policy and value function, capturing such spurious correlations can result in policies that do not generalize well to new instances, similar to what happens in supervised learning (Arjovsky et al., 2019).

Furthermore, if the environment is partially observed, an agent should have no way of predicting its expected return in a new level. At the same time, the agent could still select the optimal action if it has learned good state representations (*i.e.* that capture the minimal set of features needed to act in the environment in order to solve the task). Thus, in partially observed procedurally generated environments, accurately predicting the value function can require instance-specific features which are not necessary for learning the optimal policy.

To address the policy-value representation asymmetry, we propose **Invariant Decoupled Advantage Actor-Critic** or **IDAAC** for short, which makes two algorithmic contributions. First, IDAAC decouples the policy and value optimization by using two separate networks to learn each of them. The policy network has two heads, one for the policy and one for the generalized advantage function. The value

network is needed to compute the advantage, which is used both for the policy-gradient objective and as a target for the advantage predictions. Second, IDAAC uses an auxiliary loss which constrains the policy representation to be invariant to the task instance.

To summarize, our work makes the following contributions: (i) identifies that using a shared representation for the policy and value function can lead to overfitting in RL; (ii) proposes a new approach that uses separate networks for the policy and value function while still learning effective behaviors; (iii) introduces an auxiliary loss for encouraging representations to be invariant to the task instance, and (iv) demonstrates state-of-the-art generalization on the Procgen benchmark and outperforms popular RL methods on DeepMind Control tasks with distractors.

## 2. Related Work

**Generalization in Deep RL.** A recent body of work has pointed out the problem of overfitting in deep RL (Rajeswaran et al., 2017; Machado et al., 2018; Justesen et al., 2018; Packer et al., 2018; Zhang et al., 2018a;b; Nichol et al., 2018; Cobbe et al., 2018; 2019; Julianj et al., 2019; Raileanu & Rocktäschel, 2020; Kuttler et al., 2020; Grigsby & Qi, 2020). A promising approach to prevent overfitting is to apply regularization techniques originally developed for supervised learning such as dropout (Srivastava et al., 2014; Igl et al., 2019), batch normalization (Ioffe & Szegedy, 2015; Farebrother et al., 2018; Igl et al., 2019), or data augmentation (Cobbe et al., 2018; Ye et al., 2020; Lee et al., 2020; Laskin et al., 2020; Raileanu et al., 2020; Wang et al., 2020). Other methods use representation learning techniques to improve generalization in RL (Igl et al., 2019; Sonar et al., 2020; Stooke et al., 2020). For example, Zhang et al. (2020a;c) and Agarwal et al. (2021) learn state abstractions using various bisimulation metrics, while Roy & Konidaris (2020) align the features of two domains using Wasserstein distance. Other approaches for improving gen-

eralization in RL consist of reducing the non-stationarity inherent in RL using policy distillation (Igl et al., 2020), minimizing surprise (Chen, 2020), maximizing the mutual information between the agent’s internal representation of successive time steps (Mazouze et al., 2020), or generating an automatic curriculum for replaying different levels based on the agent’s learning potential (Jiang et al., 2020). Recently, Bengio et al. (2020) study the link between generalization and interference in TD-learning, while Bertrán et al. (2020) prove that agents trained with off-policy actor-critic methods overfit to their training instances, which is consistent with our empirical results. However, none of these works focus on the asymmetry between the optimal policy and value representation.

**Decoupling the Policy and Value Function.** While the current standard practice in deep RL from pixels is to share parameters between the policy and value function in order to learn good representations and reduce computational complexity (Mnih et al., 2016; Silver et al., 2017; Schulman et al., 2017), a few papers have explored the idea of decoupling the two for improving sample efficiency (Barth-Maron et al., 2018; Pinto et al., 2018; Yarats et al., 2019; Andrychowicz et al., 2020; Cobbe et al., 2020). In contrast with prior work, our paper focuses on generalization to unseen environments and is the first one to point out that using shared features for the policy and value functions can lead to overfitting to spurious correlations. Most similar to our work, Cobbe et al. (2020) aim to alleviate the interference between policy and value optimization, but there are some key differences between their approach and ours. In particular, our method does not require an auxiliary learning phase for distilling the value function while constraining the policy, it does not use gradients from the value function to update the policy parameters, and it uses two auxiliary losses for training the policy network, one based on the advantage function and one that enforces invariance with respect to the environment instance. Prior work has also explored the idea of predicting the advantage in the context of Q-learning (Wang et al., 2016), but this setting does not pose the same challenges since it does not learn policies directly.

### 3. Background

We consider a distribution  $q(m)$  of Partially Observable Markov Decision Processes (POMDPs)  $m \in \mathcal{M}$ , with  $m$  defined by the tuple  $(\mathcal{S}_m, \mathcal{O}_m, \mathcal{A}, \mathcal{T}_m, \Omega_m \mathcal{R}_m, \gamma)$ , where  $\mathcal{S}_m$  is the state space,  $\mathcal{O}_m$  is the observation space,  $\mathcal{A}$  is the action space,  $\mathcal{T}_m(s'|s, a)$  is the state transition probability distribution,  $\Omega_m(o'|s, a)$  is the observation probability distribution,  $\mathcal{R}_m(s, a)$  is the reward function, and  $\gamma$  is the discount factor. During training, we restrict access to a fixed set of POMDPs,  $\mathcal{M}_{train} = \{m_1, \dots, m_n\}$ , where  $m_i \sim q$ ,

$\forall i = \overline{1, n}$ . The goal is to find a policy  $\pi_\theta$  which maximizes the expected discounted reward over the entire distribution of POMDPs,  $J(\pi_\theta) = \mathbb{E}_{q, \pi, \mathcal{T}_m} [\sum_{t=0}^T \gamma^t R_m(s_t, a_t)]$ .

In practice, we use the Procgen benchmark which contains 16 procedurally generated games. Each game corresponds to a distribution of POMDPs  $q(m)$ , and each level of a game corresponds to a POMDP sampled from that game’s distribution  $m \sim q$ . The POMDP  $m$  is determined by the seed (*i.e.* integer) used to generate the corresponding level. Following the setup from Cobbe et al. (2019), agents are trained on a fixed set of  $n = 200$  levels (generated using seeds from 1 to 200) and tested on the full distribution of levels (generated using any computer integer seed).

## 4. Invariant Decoupled Advantage Actor-Critic

We start by describing our first contribution, namely the Decoupled Advantage Actor-Critic (DAAC) algorithm, which uses separate networks for learning the policy and value function (Figure 2 (left), Section 4.2). Then, we extend this method by adding an auxiliary loss to constrain the policy representation to be invariant to the environment instance, which yields the Invariant Decoupled Advantage Actor-Critic (IDAAC) algorithm (Figure 2 (right), Section 4.3).

### 4.1. Decoupling the Policy and Value Function

To address the problem of overfitting due to the coupling of the policy and value function, we propose a new actor-critic algorithm that uses separate networks for learning the policy and value function, as well as two auxiliary losses. A naive solution to the problem of parameter sharing between the policy and value function would be to simply train two separate networks. However, this approach is insufficient for learning effective behaviors because the policy network relies on gradients from the value function to learn useful features for the training environments. As shown by Cobbe et al. (2020), using separate networks for optimizing the policy and value function leads to drastically worse training performance than using a shared network, with no sign of progress on many of the tasks (see Figure 8 in their paper). Since such approaches cannot even learn useful behaviors for the training environments, they are no better on the test environments. These results indicate that without gradients from the value to update the policy network, the agent struggles to learn good behaviors. This is consistent with the fact that the gradients from the policy objective are notoriously sparse and high-variance, making training difficult, especially for high-dimensional state spaces (*e.g.* when learning from images). In contrast, the value loss can provide denser and less noisy gradients, leading to more efficient training. Given this observation, it is natural to investigate whether

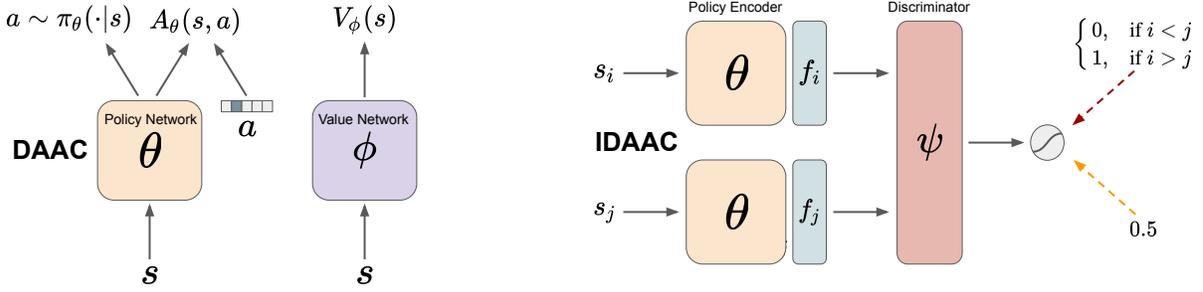


Figure 2. **Overview of DAAC (left) and IDAAC (right)**. DAAC uses two separate networks, one for learning the policy and advantage, and one for learning the value. The value estimates are used to compute the advantage targets. IDAAC adds an additional regularizer to the DAAC policy encoder to ensure that it does not contain episode-specific information. The encoder is trained adversarially with a discriminator so that it cannot classify which observation from a given pair  $(s_i, s_j)$  was first in a trajectory.

other auxiliary losses can provide similarly useful gradients for training the policy network, while also alleviating the problem of overfitting to spurious correlations.

#### 4.2. Using the Advantage instead of the Value Function

As an alternative to the value function, we propose to predict the generalized advantage estimate (GAE) or advantage for short. As illustrated in Section 5.5 and Appendix G, the advantage is less prone to overfitting to certain types of environment idiosyncrasies. Intuitively, the advantage is a measure of the expected additional return which can be obtained by taking a particular action relative to following the current policy. Because the advantage is a *relative* measure of an action’s value while the value is an *absolute* measure of a state’s value, the advantage can be expected to vary less with the number of remaining steps in the episode. Thus, the advantage is less likely to overfit to such instance-specific features. To learn generalizable representations, we need to fit a metric invariant to cosmetic changes in the observation which do not modify the underlying state. As shown in Appendix G, semantically identical yet visually distinct observations can have very different values but the same advantages (for a given action). This indicates that the advantage might be a good candidate to replace the value as an auxiliary loss for training the policy network.

In order to predict the advantage, we need an estimate of the value function which we obtain by simply training a separate network to output the expected return for a given state. Thus, our method consists of two separate networks, the value network parameterized by  $\phi$  which is trained to predict the value function, and the policy network parameterized by  $\theta$  which is trained to learn a policy that maximizes the expected return and also to predict the advantage function.

The policy network of DAAC is trained to maximize the

following objective:

$$J_{\text{DAAC}}(\theta) = J_{\pi}(\theta) + \alpha_s S_{\pi}(\theta) - \alpha_a L_A(\theta), \quad (1)$$

where  $J_{\pi}(\theta)$  is the policy gradient objective,  $S_{\pi}(\theta)$  is an entropy bonus to encourage exploration,  $L_A(\theta)$  is the advantage loss, while  $\alpha_s$  and  $\alpha_a$  are their corresponding weights determining each term’s contribution to the total objective.

The policy objective term is the same as the one used by PPO (see Appendix A):

$$J_{\pi}(\theta) = \hat{\mathbb{E}}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip} \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right],$$

where  $r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$  and  $\hat{A}_t$  is the advantage function at time step  $t$ .

The advantage function loss term is defined as:

$$L_A(\theta) = \hat{\mathbb{E}}_t \left[ \left( A_{\theta}(s_t, a_t) - \hat{A}_t \right)^2 \right],$$

where  $\hat{A}_t$  is the corresponding generalized advantage estimate at time step  $t$ ,  $\hat{A}_t = \sum_{k=t}^T (\gamma\lambda)^{k-t} \delta_k$ , with  $\delta_t = r_t + \gamma V_{\phi}(s_{t+1}) - V_{\phi}(s_t)$  which is computed using the estimates from the value network.

The value network of DAAC is trained to minimize the following loss:

$$L_V(\phi) = \hat{\mathbb{E}}_t \left[ \left( V_{\phi}(s_t) - \hat{V}_t \right)^2 \right],$$

where  $\hat{V}_t$  is the total discounted reward obtained during the corresponding episode after step  $t$ ,  $\hat{V}_t = \sum_{k=t}^T \gamma^{k-t} r_k$ .

During training, we alternate between  $E_{\pi}$  epochs for training the policy network and  $E_V$  epochs for training the value network every  $N_{\pi}$  policy updates. See Algorithm 1 from Appendix B for a more detailed description of DAAC.

As our experiments show, predicting the advantage rather than the value provides useful gradients for the policy network so it can learn effective behaviors on the training environments, thus overcoming the challenges encountered by prior attempts at learning separate policy and value networks (Cobbe et al., 2020). In addition, it mitigates the problem of overfitting caused by the use of value gradients to update the policy network, thus also achieving better performance on test environments.

### 4.3. Learning Instance-Invariant Features

From a generalization perspective, a good state representation is characterized by one that captures the minimum set of features necessary to learn the optimal policy and ignores instance-specific features which might lead to overfitting. As emphasized in Figure 1, due to the diversity of procedurally generated environments, the observations may contain information indicative of the number of remaining steps in the corresponding level. Since different levels have different lengths, capturing such information given only a partial view of the environment translates into capturing information specific to that level. Because such features overfit to the idiosyncrasies of the training environments, they can result in suboptimal policies on unseen instances of the same task.

Hence, one way of constraining the learned representations to be agnostic to the environment instance is to discourage them from carrying information about the number of remaining steps in the level. This can be formalized using an adversarial framework so that a discriminator cannot tell which observation from a given pair came first within an episode, based solely on their learned features. Similar ideas have been proposed for learning disentangled representations of videos (Denton & Birodkar, 2017).

Let  $E_\theta$  be an encoder that takes as input an observation  $s$  and outputs a feature vector  $f$ . This encoder is the same as the one used by the policy network so it is also parameterized by  $\theta$ . Let  $D$  be a discriminator parameterized by  $\psi$  that takes as input two features  $f_i$  and  $f_j$  (in this order), corresponding to two observations from the same trajectory  $s_i$  and  $s_j$ , and outputs a number between 0 and 1 which represents the probability that observation  $s_i$  came before observation  $s_j$ . The discriminator is trained using a cross-entropy loss that aims to predict which observation was first in the trajectory:

$$L_D(\psi) = -\log [D_\psi (E_\theta(s_i), E_\theta(s_j))] - \log [1 - D_\psi (E_\theta(s_i), E_\theta(s_j))]. \quad (2)$$

Note that only the discriminator’s parameters are updated by minimizing the loss in eq. 2, while the encoder’s parameters remain fixed during this optimization.

The other half of the adversarial framework imposes a loss

function on the encoder that tries to maximize the uncertainty (*i.e.* entropy) of the discriminator regarding which observation was first in the episode:

$$L_E(\theta) = -\frac{1}{2} \log [D_\psi (E_\theta(s_i), E_\theta(s_j))] - \frac{1}{2} \log [1 - D_\psi (E_\theta(s_i), E_\theta(s_j))]. \quad (3)$$

Similar to the above, only the encoder’s parameters are updated by minimizing the loss in eq. 3, while the discriminator’s parameters remain fixed during this optimization.

Thus, the policy network is encouraged to learn state representations so that the discriminator cannot identify whether a state came before or after another state. In so doing, the learned representations cannot carry information about the number of remaining steps in the environment, yielding features which are less instance-dependent and thus more likely to generalize outside the training distribution. Note that this adversarial loss is only used for training the policy network and not the value network.

To train the policy network, we maximize the following objective which combines the DAAC objective from eq. 1 with the above adversarial loss, resulting in IDAAC’s objective:

$$J_{IDAAC}(\theta) = J_{DAAC}(\theta) - \alpha_i L_E(\theta), \quad (4)$$

where  $\alpha_i$  is the weight of the adversarial loss relative to the policy objective. Similar to DAAC, a separate value network is trained. See Algorithm 2 from Appendix B for a more detailed description of IDAAC.

## 5. Experiments

In this section, we evaluate our methods on two distinct environments: (i) three DeepMind Control suite tasks with synthetic and natural background distractors (Zhang et al., 2020b) and (ii) the full Procgen benchmark (Cobbe et al., 2019) which consists of 16 procedurally generated games. Procgen in particular has a number of attributes that make it a good testbed for generalization in RL: (i) it has a diverse set of games in a similar spirit with the ALE benchmark (Bellemare et al., 2013); (ii) each of these games has procedurally generated levels which present agents with meaningful generalization challenges; (iii) agents have to learn motor control directly from images, and (iv) it has a clear protocol for testing generalization, the focus of our investigation.

All Procgen environments use a discrete 15 dimensional action space and produce  $64 \times 64 \times 3$  RGB observations. We use Procgen’s *easy* setup, so for each game, agents are trained on 200 levels and tested on the full distribution of levels. More details about our experimental setup and hyperparameters can be found in Appendix C.

Table 1. PPO-Normalized Procgen scores on train and test levels after training on 25M environment steps. Our approaches, DAAC and IDAAC, establish a new state-of-the-art on the test distribution of environments from the Procgen benchmark, while also showing strong training performance. The mean and standard deviation are computed using 10 runs with different seeds.

Score	RAND-FM	IBAC-SNI	Mixreg	PLR	UCB-DrAC	PPG	DAAC (Ours)	IDAAC (Ours)
Train	87.6 ± 8.9	103.4 ± 8.5	104.2 ± 3.1	106.7 ± 5.6	118.9 ± 8.0	<b>144.5 ± 5.7</b>	131.0 ± 6.1	132.2 ± 5.9
Test	78.0 ± 9.0	102.9 ± 8.6	114.6 ± 3.3	128.3 ± 5.8	139.7 ± 8.3	152.2 ± 5.8	162.3 ± 6.2	<b>163.7 ± 6.1</b>

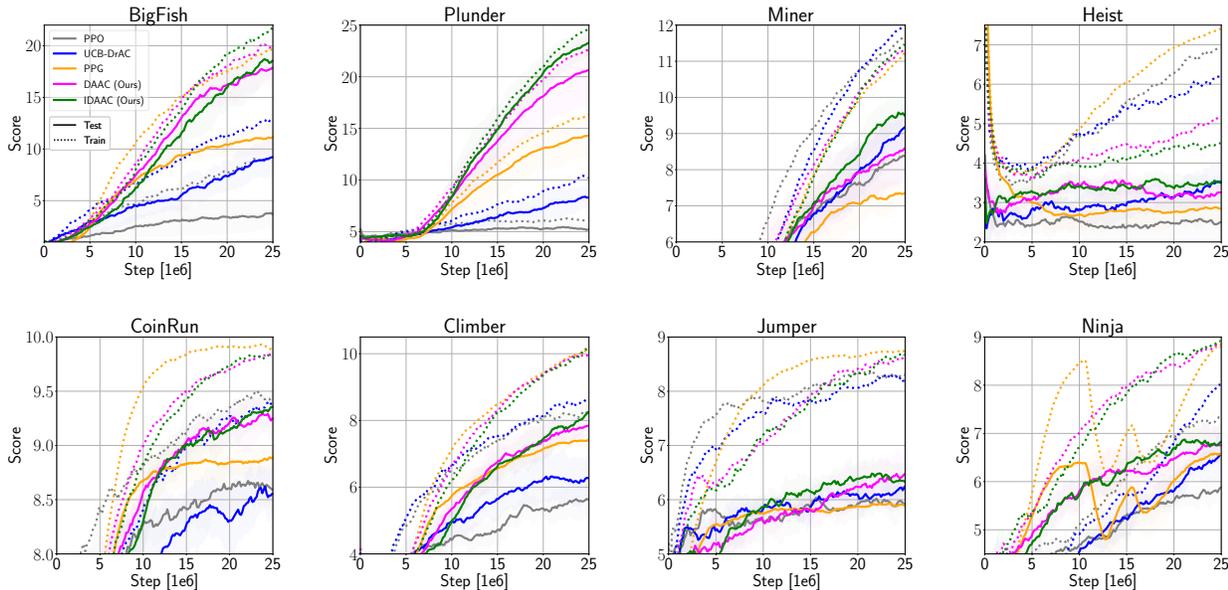


Figure 3. Train and Test Performance for IDAAC, DAAC, PPG, UCB-DrAC, and PPO, on eight diverse Procgen games. IDAAC and DAAC display state-of-the-art performance on the test levels, beating leading approaches (PPG and UCB-DrAC), as well a PPO baseline. Furthermore, IDAAC and DAAC exhibit a smaller generalization gap than other methods. The mean and standard deviation are computed over 10 runs with different seeds.

### 5.1. Generalization Performance on Procgen

We compare DAAC and IDAAC with seven other RL algorithms: **PPO** (Schulman et al., 2017), **UCB-DrAC** (Raileanu et al., 2020), **PLR** (Jiang et al., 2020), **Mixreg** (Wang et al., 2020), **IBAC-SNI** (Igl et al., 2019), **Rand-FM** (Lee et al., 2020), and **PPG** (Cobbe et al., 2020).

UCB-DrAC is the previous state-of-the-art on Procgen and uses data augmentation to learn policy and value functions invariant to various input transformations PLR is a newer approach that uses an automatic curriculum based on the learning potential of each level and achieves strong results on Procgen. Rand-FM uses a random convolutional network to regularize the learned representations, IBAC-SNI uses an information bottleneck with selective noise injection, while Mixreg uses mixtures of observations to impose linearity constraints between the agent’s inputs and outputs. All three were designed to improve generalization in RL and evaluated on Procgen games. Finally, PPG is the only method we are aware of that learns good policies while decoupling

the optimization of the policy and value function. However, PPG was designed to improve sample efficiency rather than generalization and the method was evaluated only on Procgen’s training distribution of environments. See Section 2 for a detailed discussion of the differences between PPG and our methods.

Table 1 shows the train and test performance of all methods, aggregated across all Procgen games. DAAC and IDAAC outperform all the baselines on both train and test. Figure 3 shows the train and test performance on eight of the Procgen games. We show comparisons with a vanilla RL algorithm PPO, the previous state-of-the-art UCB-DrAC and our strongest baseline PPG. Both of our approaches, DAAC and IDAAC, show superior results on the test levels, relative to the other methods. In addition, IDAAC and DAAC achieve better or comparable performance on the training levels for most of these games. While DAAC already shows notable gains over the baselines, IDAAC further improves upon it, thus emphasizing the benefits of combining our two contributions. See Appendix D for results on all games.

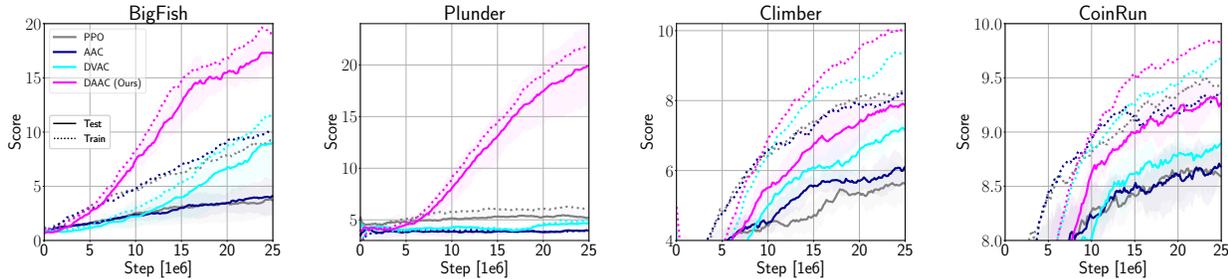


Figure 4. Train and Test Performance for PPO, DAAC, and two of its ablations, DVAC and AAC, on four Procgen games. DAAC outperforms all the ablations on both train and test, emphasizing the importance of each component. The mean and standard deviation are computed over 5 runs with different seeds.

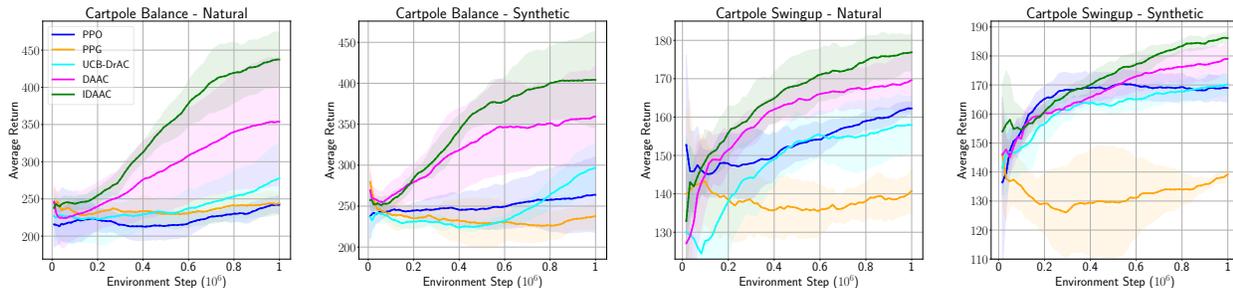


Figure 5. Average return on two DMC tasks, Cartpole Balance and Cartpole Swingup with natural and synthetic video backgrounds. Our DAAC and IDAAC approaches outperform PPO, PPG, and UCB-DrAC. The mean and standard deviation are computed over 10 runs with different seeds.

### 5.2. Ablations

We also performed a number of ablations to emphasize the importance of each component used by our method. First, **Decoupled Value Actor-Critic** or **DVAC** is an ablation to **DAAC** that learns to predict the value rather than the advantage for training the policy network. This ablation helps disentangle the effect of predicting the advantage function from the effect of using a separate value network and performing multiple updates for the value than for the policy. In principle, this decoupling could result in a more accurate estimate of the value function, which can in turn lead to more effective policy optimization. Second, **Advantage Actor-Critic** or **AAC** is a modification to PPO that includes an extra head for predicting the advantage function (in a single network). The role of this ablation is to understand the importance of *not* backpropagating gradients from the value into the policy network, even while using gradients from the advantage.

Figure 4 shows the train and test performance of DAAC, DVAC, AAC, and PPO on four Procgen games. DAAC outperforms all the ablations on both train and test environments, emphasizing the importance of each component. In particular, the fact that AAC’s generalization ability is worse

than that of DAAC suggests that predicting the advantage function in addition to also predicting the value function (as part of training the policy network) does not solve the problem of overfitting. Hence, using the value loss to update the policy parameters still hurts generalization even when the advantage is also used to train the network. In addition, the fact that DVAC has worse test performance than DAAC indicates that DAAC’s gains are not merely due to having access to a more accurate value function or to the reduced interference between optimizing the policy and value.

These results are consistent with our claim that using gradients from the value to update the policy can lead to representations that overfit to spurious correlations in the training environments. In contrast, using gradients from the advantage to train the policy network leads to agents that generalize better to new environments.

### 5.3. DeepMind Control with Distractors

In this section, we evaluate our methods on the DeepMind Control Suite from pixels (DMC, Tassa et al. (2018)). We use three tasks, namely Cartpole Balance, Cartpole Swingup, and Ball In Cup. For each task, we study two settings with different types of backgrounds, namely *synthetic* distractors

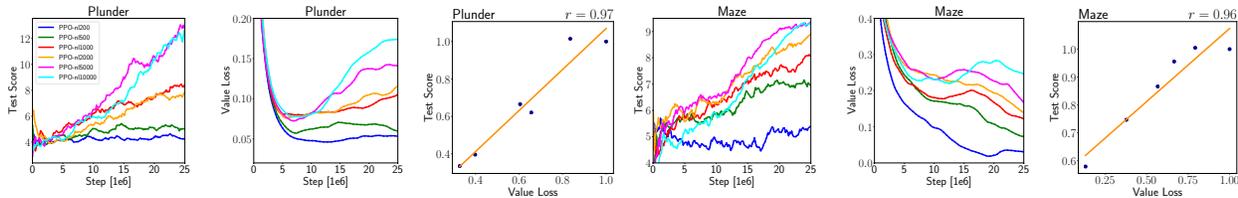


Figure 6. PPO agents trained on varying numbers of levels: 200, 500, 1000, 2000, 5000, and 10000 for Plunder (left) and Maze (right). For each game, from left to right we show the test score, value loss, and correlation between value loss and test score after 25M training steps. Both the test score and the value loss increase with the number of training levels used. These results indicate there is a positive correlation between value loss and generalization when using a shared network for the policy and value function.

and natural videos from the Kinetics dataset (Kay et al., 2017), as introduced in Zhang et al. (2020b). Note that in the synthetic and natural settings, the background is sampled from a list of videos at the beginning of each episode, which creates spurious correlations between backgrounds and rewards. As shown in Figure 5, DAAC and IDAAC significantly outperform PPO, UCB-DrAC, and PPG on all these environments. See Appendix E for more details about the experimental setup and results on other DMC tasks.

#### 5.4. Value Loss and Generalization

When using actor-critic policy-gradient algorithms, a more accurate estimate of the value function leads to lower variance gradients and thus better policy optimization on a given environment (Sutton et al., 1999). While an accurate value function improves sample efficiency and training performance, it can also lead to overfitting when the policy and value share the same representation. To validate this claim, we looked at the correlation between value loss and test performance. More specifically, we trained 6 PPO agents on varying numbers of Procgen levels, namely 200, 500, 1000, 2000, 5000, and 10000. As expected, models trained on a larger number of levels generalize better to unseen levels as illustrated in Figure 6. However, *agents trained on more levels have a higher value loss at the end of training than agents trained on fewer levels, so the value loss is positively correlated with generalization ability*. This observation is consistent with our claim that using a shared network for the policy and value function can lead to overfitting. Our hypothesis was that, when using a common network for the policy and value function, accurately predicting the values implies that the learned representation relies on spurious correlation, which would likely lead to poor generalization at test time. Similarly, an agent with good generalization suggests that its representation relies on the features needed to learn an optimal policy for the entire family of environments, which are insufficient for accurately predicting the value function (as explained in Section 1.1). See Appendix F for the relationship between value loss, test score, and the number of training levels for all Procgen games. In Ap-

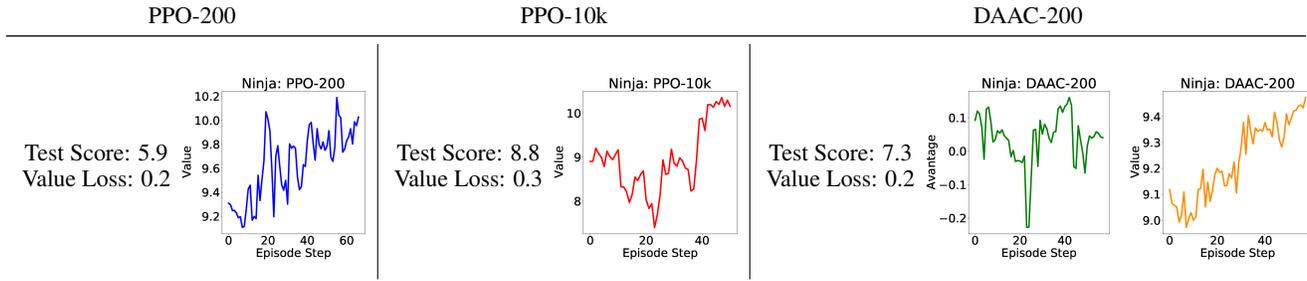
pendix H, you can see that the variance of predicted values for the initial observations decreases with the number of training levels, which further supports our claim.

#### 5.5. Advantage vs. Value During an Episode

In Section 1, we claim that in procedurally generated environments with partial observability, in order to accurately estimate the value function, the agent needs to memorize the number of remaining steps in the level. As Figure 1 shows, a standard RL agent predicts very different values for the initial states of two levels even if the observations are semantically identical. This suggests that the agent must have memorized the length of each level since the partial observation at the beginning of a level does not contain enough information to accurately predict the expected return.

We further investigate this issue by exploring the predicted value’s dependency on the episode step. Instead of just comparing the initial states as in Figure 1, we plot the value predicted by the agent over the course of an entire trajectory in one of the Ninja levels (see Table 2). Since the agent is rewarded only if it reaches the goal at the end of the episode, the true value increases linearly with the episode step over the course of the agent’s trajectory. As seen in Table 2, the estimated value of a PPO agent trained on 200 levels also has a quasi-linear dependence on the episode step, suggesting that the agent must know how many remaining steps are in the game at any point during this episode. However, the episode step cannot be inferred solely from partial observations since the training levels contain observations with the same semantics but different values, as illustrated in Figure 1. In order to accurately predict the values of such observations, the agent must learn representations that capture level-specific features (such as the backgrounds) which would allow it to differentiate between semantically similar observations with different values. Since PPO uses a common representation for the policy and value function, this can lead to policies that overfit to the particularities of the training environments. Note that a PPO agent trained on 10k levels does not show the same linear trend between the value and episode step. This implies that there is a trade-off

Table 2. The trade-off between generalization and value accuracy illustrated on a single Ninja level. A PPO agent trained on 200 levels (blue) has high value accuracy but low generalization performance, and its value predictions have a near linear dependency on the episode step. This linear relationship further supports the claim that the agent memorizes level-specific features which are needed to predict the value function given only partial observations. In contrast, a PPO agent trained on 10k levels (red) with good generalization but low value accuracy does not display this linear trend. When sharing parameters for the policy and value function, there is a trade-off between fitting the value and learning general policies. By decoupling the policy and value, our model DAAC can achieve both high value accuracy and good generalization performance. To train the policy, DAAC uses gradients from predicting the advantage (green), which does not display the linear trend, thus is less prone to overfitting. DAAC’s value estimate (orange) still shows a linear trend but, in contrast to PPO, this does not negatively affect the policy since DAAC uses separate networks to learn the policy and value.



between generalization and value accuracy for models that use a shared network to learn the policy and value.

By decoupling the policy and value, our model DAAC can achieve both high value accuracy and good generalization performance. Note that the advantage estimated by DAAC (trained on 200 levels) shows no clear dependence on the environment step, thus being less prone to overfitting. Nevertheless, DAAC’s value estimate still shows a linear trend but, in contrast to PPO, this does not negatively affect the policy since DAAC uses separate networks for the policy and value. This analysis indicates that using advantages rather than values to update the policy network leads to better generalization performance while also being able to accurately predict the value function. See Appendix G for similar results on other Procgen games, as well as comparisons with PPG which displays a similar trend as PPO trained on 200 levels.

## 6. Discussion

In this work, we identified a new problem with standard deep reinforcement learning algorithms which causes overfitting, namely the asymmetry between the policy and value representation. To alleviate this problem, we propose IDAAC, which decouples the optimization of the policy and value function while still learning effective behaviors. IDAAC also introduces an auxiliary loss which constrains the policy representation to be invariant with respect to the environment instance. IDAAC achieves a new state-of-the-art on the Procgen benchmark and outperforms strong RL algorithms on DeepMind Control tasks with distractors. In contrast to other popular methods, our approach can both achieve good generalization while also learning accurate value estimates. Moreover, IDAAC learns representations

and predictions which are more robust to cosmetic changes in the observations that do not change the underlying state of the environment (see Appendix I).

One limitation of our work is the focus on learning representations which are invariant to the number of remaining steps in the episode. While this inductive bias will not be helpful for all problems, the settings where we can expect most gains are those with partial observability, a set of goal states, and episode length variations (e.g. navigation of different layouts). A promising avenue for future work is to investigate other auxiliary losses in order to efficiently learn more general behaviors. One desirable property of such auxiliary losses is to capture the minimal set of features needed to act in the environment. While our experiments show that predicting the advantage function improves generalization, we currently lack a firm theoretical argument for this. The advantage could act as a regularizer, being less prone to memorizing the remaining episode length, or it could be better correlated with the underlying state of the environment rather than its visual appearance. Investigating these hypotheses could further improve our understanding of what leads to better representations and what we are still missing. Finally, the solution we propose here is only a first step towards solving the policy-value representation asymmetry and we hope many other ideas will be explored in future work.

## Acknowledgements

We would like to thank our ICML reviewers, as well as Vitaly Kurin, Denis Yarats, Mahi Shafiullah, Ilya Kostrikov, and Max Goldstein for their valuable feedback on this work. Roberta was supported by the DARPA Machine Commonsense program.

## References

- Agarwal, R., Machado, M. C., Castro, P. S., and Bellemare, M. G. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. 2021.
- Andrychowicz, M., Raichuk, A., Stanczyk, P., Orsini, M., Girgin, S., Marinier, R., Hussenot, L., Geist, M., Pietquin, O., Michalski, M., Gelly, S., and Bachem, O. What matters in on-policy reinforcement learning? a large-scale empirical study. *ArXiv*, abs/2006.05990, 2020.
- Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. Invariant risk minimization. *ArXiv*, abs/1907.02893, 2019.
- Barth-Maron, G., Hoffman, M. W., Budden, D., Dabney, W., Horgan, D., Dhruva, T., Muldal, A., Heess, N., and Lillicrap, T. Distributed distributional deterministic policy gradients. *ArXiv*, abs/1804.08617, 2018.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- Bengio, E., Pineau, J., and Precup, D. Interference and generalization in temporal difference learning. *ArXiv*, abs/2003.06350, 2020.
- Bertrán, M., Martínez, N., Phielipp, M., and Sapiro, G. Instance based generalization in reinforcement learning. *ArXiv*, abs/2011.01089, 2020.
- Chen, J. Z. Reinforcement learning generalization with surprise minimization. *ArXiv*, abs/2004.12399, 2020.
- Cobbe, K., Klimov, O., Hesse, C., Kim, T., and Schulman, J. Quantifying generalization in reinforcement learning. *arXiv preprint arXiv:1812.02341*, 2018.
- Cobbe, K., Hesse, C., Hilton, J., and Schulman, J. Leveraging procedural generation to benchmark reinforcement learning. *arXiv preprint arXiv:1912.01588*, 2019.
- Cobbe, K., Hilton, J., Klimov, O., and Schulman, J. Phasic policy gradient. *ArXiv*, abs/2009.04416, 2020.
- Denton, E. L. and Birodkar, V. Unsupervised learning of disentangled representations from video. In *NIPS*, 2017.
- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. *arXiv preprint arXiv:1802.01561*, 2018.
- Farebrother, J., Machado, M. C., and Bowling, M. H. Generalization and regularization in dqn. *ArXiv*, abs/1810.00123, 2018.
- Gamrian, S. and Goldberg, Y. Transfer learning for related reinforcement learning tasks via image-to-image translation. *ArXiv*, abs/1806.07377, 2019.
- Grigsby, J. and Qi, Y. Measuring visual generalization in continuous control from pixels. *ArXiv*, abs/2010.06740, 2020.
- Igl, M., Ciosek, K., Li, Y., Tschitschek, S., Zhang, C., Devlin, S., and Hofmann, K. Generalization in reinforcement learning with selective noise injection and information bottleneck. In *Advances in Neural Information Processing Systems*, pp. 13956–13968, 2019.
- Igl, M., Farquhar, G., Luketina, J., Böhrer, W., and Whiteson, S. The impact of non-stationarity on generalisation in deep reinforcement learning. *ArXiv*, abs/2006.05826, 2020.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ArXiv*, abs/1502.03167, 2015.
- Jiang, M., Grefenstette, E., and Rocktäschel, T. Prioritized level replay. *ArXiv*, abs/2010.03934, 2020.
- Juliani, A., Khalifa, A., Berges, V.-P., Harper, J., Henry, H., Crespi, A., Togelius, J., and Lange, D. Obstacle tower: A generalization challenge in vision, control, and planning. *ArXiv*, abs/1902.01378, 2019.
- Justesen, N., Torrado, R. R., Bontrager, P., Khalifa, A., Togelius, J., and Risi, S. Illuminating generalization in deep reinforcement learning through procedural level generation. *arXiv: Learning*, 2018.
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, A., Suleyman, M., and Zisserman, A. The kinetics human action video dataset. *ArXiv*, abs/1705.06950, 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.
- Kostrikov, I. Pytorch implementations of reinforcement learning algorithms. <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail>, 2018.
- Kostrikov, I., Yarats, D., and Fergus, R. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.
- Kuttler, H., Nardelli, N., Miller, A. H., Raileanu, R., Selvatici, M., Grefenstette, E., and Rocktäschel, T. The nethack learning environment. *ArXiv*, abs/2006.13760, 2020.

- Laskin, M., Lee, K., Stooke, A., Pinto, L., Abbeel, P., and Srinivas, A. Reinforcement learning with augmented data. *arXiv preprint arXiv:2004.14990*, 2020.
- Lee, K., Lee, K., Shin, J., and Lee, H. Network randomization: A simple technique for generalization in deep reinforcement learning. In *International Conference on Learning Representations*. <https://openreview.net/forum/2020>.
- Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M. J., and Bowling, M. H. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. In *IJCAI*, 2018.
- Mazouze, B., des Combes, R. T., Doan, T., Bachman, P., and Hjelm, R. D. Deep reinforcement and infomax learning. *ArXiv*, abs/2006.07217, 2020.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. *ArXiv*, abs/1602.01783, 2016.
- Nichol, A., Pfau, V., Hesse, C., Klimov, O., and Schulman, J. Gotta learn fast: A new benchmark for generalization in rl. *ArXiv*, abs/1804.03720, 2018.
- Packer, C., Gao, K., Kos, J., Krähenbühl, P., Koltun, V., and Song, D. X. Assessing generalization in deep reinforcement learning. *ArXiv*, abs/1810.12282, 2018.
- Pinto, L., Andrychowicz, M., Welinder, P., Zaremba, W., and Abbeel, P. Asymmetric actor critic for image-based robot learning. *ArXiv*, abs/1710.06542, 2018.
- Raileanu, R. and Rocktäschel, T. Ride: Rewarding impact-driven exploration for procedurally-generated environments. *ArXiv*, abs/2002.12292, 2020.
- Raileanu, R., Goldstein, M., Yarats, D., Kostrikov, I., and Fergus, R. Automatic data augmentation for generalization in deep reinforcement learning. *ArXiv*, abs/2006.12862, 2020.
- Rajeswaran, A., Lowrey, K., Todorov, E., and Kakade, S. M. Towards generalization and simplicity in continuous control. *ArXiv*, abs/1703.02660, 2017.
- Roy, J. and Konidaris, G. Visual transfer for reinforcement learning via wasserstein domain confusion. *arXiv preprint arXiv:2006.03465*, 2020.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M. I., and Moritz, P. Trust region policy optimization. In *ICML*, 2015.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., and Hassabis, D. Mastering the game of go without human knowledge. *Nature*, 550:354–359, 2017.
- Sonar, A., Pacelli, V., and Majumdar, A. Invariant policy optimization: Towards stronger generalization in reinforcement learning. *ArXiv*, abs/2006.01096, 2020.
- Song, X., Jiang, Y., Tu, S., Du, Y., and Neyshabur, B. Observational overfitting in reinforcement learning. *ArXiv*, abs/1912.02975, 2020.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15:1929–1958, 2014.
- Stooke, A., Lee, K., Abbeel, P., and Laskin, M. Decoupling representation learning from reinforcement learning. *ArXiv*, abs/2009.08319, 2020.
- Sutton, R., McAllester, D. A., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, 1999.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., Lillicrap, T., and Riedmiller, M. A. Deepmind control suite. *ArXiv*, abs/1801.00690, 2018.
- Wang, K., Kang, B., Shao, J., and Feng, J. Improving generalization in reinforcement learning with mixture regularization. *ArXiv*, abs/2010.10814, 2020.
- Wang, Z., Schaul, T., Hessel, M., Hasselt, H. V., Lanctot, M., and Freitas, N. D. Dueling network architectures for deep reinforcement learning. *ArXiv*, abs/1511.06581, 2016.
- Yarats, D., Zhang, A., Kostrikov, I., Amos, B., Pineau, J., and Fergus, R. Improving sample efficiency in model-free reinforcement learning from images. *ArXiv*, abs/1910.01741, 2019.
- Ye, C., Khalifa, A., Bontrager, P., and Togelius, J. Rotation, translation, and cropping for zero-shot generalization. *arXiv preprint arXiv:2001.09908*, 2020.
- Zhang, A., Ballas, N., and Pineau, J. A dissection of overfitting and generalization in continuous reinforcement learning. *ArXiv*, abs/1806.07937, 2018a.

Zhang, A., Lyle, C., Sodhani, S., Filos, A., Kwiatkowska, M., Pineau, J., Gal, Y., and Precup, D. Invariant causal prediction for block mdps. *arXiv preprint arXiv:2003.06016*, 2020a.

Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning invariant representations for reinforcement learning without reconstruction. 2020b.

Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020c.

Zhang, C., Vinyals, O., Munos, R., and Bengio, S. A study on overfitting in deep reinforcement learning. *ArXiv*, abs/1804.06893, 2018b.

## A. PPO

**Proximal Policy Optimization (PPO)** (Schulman et al., 2017) is an actor-critic algorithm that learns a policy  $\pi_\theta$  and a value function  $V_\theta$  with the goal of finding an optimal policy for a given MDP. PPO alternates between sampling data through interaction with the environment and optimizing an objective function using stochastic gradient ascent. At each iteration, PPO maximizes the following objective:

$$J_{\text{PPO}} = J_\pi - \alpha_1 J_V + \alpha_2 S_{\pi_\theta}, \quad (5)$$

where  $\alpha_1, \alpha_2$  are weights for the different loss terms,  $S_{\pi_\theta}$  is the entropy bonus for aiding exploration,  $J_V$  is the value function loss defined as

$$J_V = (V_\theta(s) - V_t^{\text{target}})^2.$$

The policy objective term  $J_\pi$  is based on the policy gradient objective which can be estimated using importance sampling in off-policy settings (*i.e.* when the policy used for collecting data is different from the policy we want to optimize):

$$J_{PG}(\theta) = \sum_{a \in \mathcal{A}} \pi_\theta(a|s) \hat{A}_{\theta_{\text{old}}}(s, a) = \mathbb{E}_{a \sim \pi_{\theta_{\text{old}}}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} \hat{A}_{\theta_{\text{old}}}(s, a) \right], \quad (6)$$

where  $\hat{A}(\cdot)$  is an estimate of the advantage function,  $\theta_{\text{old}}$  are the policy parameters before the update,  $\pi_{\theta_{\text{old}}}$  is the behavior policy used to collect trajectories (*i.e.* that generates the training distribution of states and actions), and  $\pi_\theta$  is the policy we want to optimize (*i.e.* that generates the true distribution of states and actions).

This objective can also be written as

$$J_{PG}(\theta) = \mathbb{E}_{a \sim \pi_{\theta_{\text{old}}}} \left[ r(\theta) \hat{A}_{\theta_{\text{old}}}(s, a) \right], \quad (7)$$

where

$$r_\theta = \frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}$$

is the importance weight for estimating the advantage function.

PPO is inspired by TRPO (Schulman et al., 2015), which constrains the update so that the policy does not change too much in one step. This significantly improves training stability and leads to better results than vanilla policy gradient algorithms. TRPO achieves this by minimizing the KL divergence between the old (*i.e.* before an update) and the new (*i.e.* after an update) policy. PPO implements the constraint in a simpler way by using a clipped surrogate objective instead of the more complicated TRPO objective. More specifically, PPO imposes the constraint by forcing  $r(\theta)$  to stay within a small interval around 1, precisely  $[1 - \epsilon, 1 + \epsilon]$ , where  $\epsilon$  is a hyperparameter. The policy objective term from equation (5) becomes

$$J_\pi = \mathbb{E}_\pi \left[ \min \left( r_\theta \hat{A}, \text{clip}(r_\theta, 1 - \epsilon, 1 + \epsilon) \hat{A} \right) \right],$$

where  $\hat{A} = \hat{A}_{\theta_{\text{old}}}(s, a)$  for brevity. The function  $\text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon)$  clips the ratio to be no more than  $1 + \epsilon$  and no less than  $1 - \epsilon$ . The objective function of PPO takes the minimum one between the original value and the clipped version so that agents are discouraged from increasing the policy update to extremes for better rewards.

## B. DAAC and IDAAC

As described in the paper, DAAC and IDAAC alternate between optimizing the policy network and optimizing the value network. The value estimates are used to compute the advantage targets which are needed by both the policy gradient objective and the auxiliary loss based on predicting the advantage function. Since we use separate networks for learning the policy and value function, we can now use different numbers of epochs for updating the two networks. We use  $E_\pi$  epochs for every policy update and  $E_V$  epochs for every value update. Similar to Cobbe et al. (2020), we find that the value network allows for larger amounts of sample reuse than the policy network. This decoupling of the policy and value optimization allows us to also control the frequency of value updates relative to policy updates. Updating the value function less often can help with training stability since can result in lower variance gradients for the policy and advantage losses. We update the value network every  $N_\pi$  updates of the policy network. See algorithms 1 and 2 for the pseudocodes of DAAC and IDAAC, respectively.

---

### Algorithm 1 DAAC: Decoupled Advantage Actor-Critic

---

```

1: Hyperparameters: Total number of updates N, replay buffer size T, number of epochs per policy update  $E_\pi$ , number of
   epochs per value update  $E_V$ , frequency of value updates  $N_\pi$ , weight for the advantage loss  $\alpha_a$ , initial policy parameters
    $\theta$ , initial value parameters  $\phi$ .
2: for  $n = 1, \dots, N$  do
3:   Collect  $\mathcal{D} = \{(s_t, a_t, r_t, s_{t+1})\}_{t=1}^T$  using  $\pi(\theta)$ .
4:   Compute the value and advantage targets  $\hat{V}_t$  and  $\hat{A}_t$  for all states  $s_t$ 
5:   for  $i = 1, \dots, E_\pi$  do
6:      $L_A(\theta) = \hat{\mathbb{E}}_t \left[ \left( A_\theta(s_t, a_t) - \hat{A}_t \right)^2 \right]$  ▷ Compute the Advantage Loss
7:      $J_{\text{DAAC}}(\theta) = J_\pi(\theta) + \alpha_s S_\pi(\theta) - \alpha_a L_A(\theta)$  ▷ Compute the Policy Loss
8:      $\theta \leftarrow \arg \max_\theta J_{\text{DAAC}}$  ▷ Update the Policy Network
9:   if  $n \% N_\pi = 0$  then
10:    for  $j = 1, \dots, E_V$  do
11:       $L_V(\phi) = \hat{\mathbb{E}}_t \left[ \left( V_\phi(s_t) - \hat{V}_t \right)^2 \right]$  ▷ Compute the Value Loss
12:       $\phi \leftarrow \arg \min_\phi L_V$  ▷ Update the Value Network

```

---

### Algorithm 2 IDAAC: Invariant Decoupled Advantage Actor-Critic

---

```

1: Hyperparameters: Total number of updates N, replay buffer size T, number of epochs per policy update  $E_\pi$ , number of
   epochs per value update  $E_V$ , frequency of value updates  $N_\pi$ , weight for the invariance loss  $\alpha_i$ , initial policy parameters
    $\theta$ , initial value parameters  $\phi$ .
2: for  $n = 1, \dots, N$  do
3:   Collect  $\mathcal{D} = \{(s_t, a_t, r_t, s_{t+1})\}_{t=1}^T$  using  $\pi(\theta)$ .
4:   Compute the value and advantage targets  $\hat{V}_t$  and  $\hat{A}_t$  for all states  $s_t$ 
5:   for  $i = 1, \dots, E_\pi$  do
6:      $L_A(\theta) = \hat{\mathbb{E}}_t \left[ \left( A_\theta(s_t, a_t) - \hat{A}_t \right)^2 \right]$  ▷ Compute the Advantage Loss
7:      $L_E(\theta) = -\frac{1}{2} \log [D_\psi(E_\theta(s_i), E_\theta(s_j))] - \frac{1}{2} \log [1 - D_\psi(E_\theta(s_i), E_\theta(s_j))]$  ▷ Compute the Encoder Loss
8:      $J_{\text{IDAAC}}(\theta, \phi, \psi) = J_\pi(\theta) + \alpha_s S_\pi(\theta) - \alpha_a L_A(\theta) - \alpha_i L_E(\theta)$  ▷ Compute the Policy Loss
9:      $L_D(\psi) = -\log [D_\psi(E_\theta(s_i), E_\theta(s_j))] - \log [1 - D_\psi(E_\theta(s_i), E_\theta(s_j))]$  ▷ Compute the Discriminator Loss
10:     $\theta \leftarrow \arg \max_\theta J_{\text{IDAAC}}$  ▷ Update the Policy Network
11:     $\psi \leftarrow \arg \min_\psi L_D$  ▷ Update the Discriminator
12:   if  $n \% N_\pi = 0$  then
13:    for  $j = 1, \dots, E_V$  do
14:       $L_V(\phi) = \hat{\mathbb{E}}_t \left[ \left( V_\phi(s_t) - \hat{V}_t \right)^2 \right]$  ▷ Compute the Value Loss
15:       $\phi \leftarrow \arg \min_\phi L_V$  ▷ Update the Value Network

```

---

Table 3. List of hyperparameters used to obtain the results in this paper.

Hyperparameter	Value
$\gamma$	0.999
$\lambda$	0.95
# timesteps per rollout	256
# epochs per rollout	3
# minibatches per epoch	8
entropy bonus	0.01
clip range	0.2
reward normalization	yes
learning rate	5e-4
# workers	1
# environments per worker	64
# total timesteps	25M
optimizer	Adam
LSTM	no
frame stack	no

### C. Hyperparameters

We use [Kostrikov \(2018\)](#)’s implementation of PPO ([Schulman et al., 2017](#)), on top of which all our methods are build. The agent is parameterized by the ResNet architecture from [Espeholt et al. \(2018\)](#) which was used to obtain the best results in [Cobbe et al. \(2019\)](#). Unless otherwise noted, we use the best hyperparameters found in [Cobbe et al. \(2019\)](#) for the easy mode of Procgen (*i.e.* same experimental setup as the one used here) as found in Table 3:

We ran a hyperparameter search over the number of epochs used during each update of the policy network  $E_\pi \in [1, 3, 6]$  the number epochs used during each update of the value network  $E_V \in [1, 5, 9]$ , the number of value updates after which we perform a policy update  $N_\pi \in [1, 8, 32]$ , the weight for the advantage loss  $\alpha_a \in [0.05, 0.1, 0.15, 0.2, 0.25, 0.3]$ , and the weight for the instance-invariant (adversarial) loss  $\alpha_i \in [0.001, 0.005, 0.01, 0.05, 0.1, 0.2]$ . We found  $E_\pi = 1, E_V = 9, N_\pi = 1, \alpha_a = 0.25, \alpha_i = 0.001$  to be the best hyperparameters overall, so we used these values to obtain the results reported in the paper. For all our experiments we use the Adam ([Kingma & Ba, 2015](#)) optimizer.

For all baselines, we use the same experimental setup for training and testing as the one used for our methods. Hence, we train them for 25M frames on the easy mode of each Procgen game, using (the same) 200 levels for training and the entire distribution of levels for testing.

For Mixreg, PLR, UCB-DrAC, and PPG, we used the best hyperparameters reported by the authors, since all these methods use Procgen for evaluation and they performed extensive hyperparameter sweeps.

For Rand-FM ([Lee et al., 2020](#)) we use the recommended hyperparameters in the authors’ released implementation, which were the best values for CoinRun ([Cobbe et al., 2018](#)), one of the Procgen games used for evaluation in ([Lee et al., 2020](#)).

For IBAC-SNI ([Igl et al., 2019](#)) we also use the authors’ open sourced implementation. We use the parameters corresponding to IBAC-SNI  $\lambda = .5$ . We use weight regularization with  $l_2 = .0001$ , data augmentation turned on, and a value of  $\beta = .0001$  which turns on the variational information bottleneck, and selective noise injection turned on. This corresponds to the best version of this approach, as found by the authors after evaluating it on CoinRun ([Cobbe et al., 2018](#)).

## D. Procgen Results

Figures 7 and 8 show the test and train performance for IDAAC, DAAC, PPG, UCB-DrAC, and PPO on all Procgen games. Our methods, IDAAC and DAAC demonstrate superior test performance, outperforming all other baselines on the majority of the games, while being comparable on most of the remaining ones.

Figures 9 and 10 show the test and train performance for PPO, DAAC, and two ablations, DVAC and AAC, on all Procgen games. On both train and test environments, our method is substantially better than all the ablations for most of the games and comparable on the remaining ones (*i.e.* CaveFlyer and Maze).

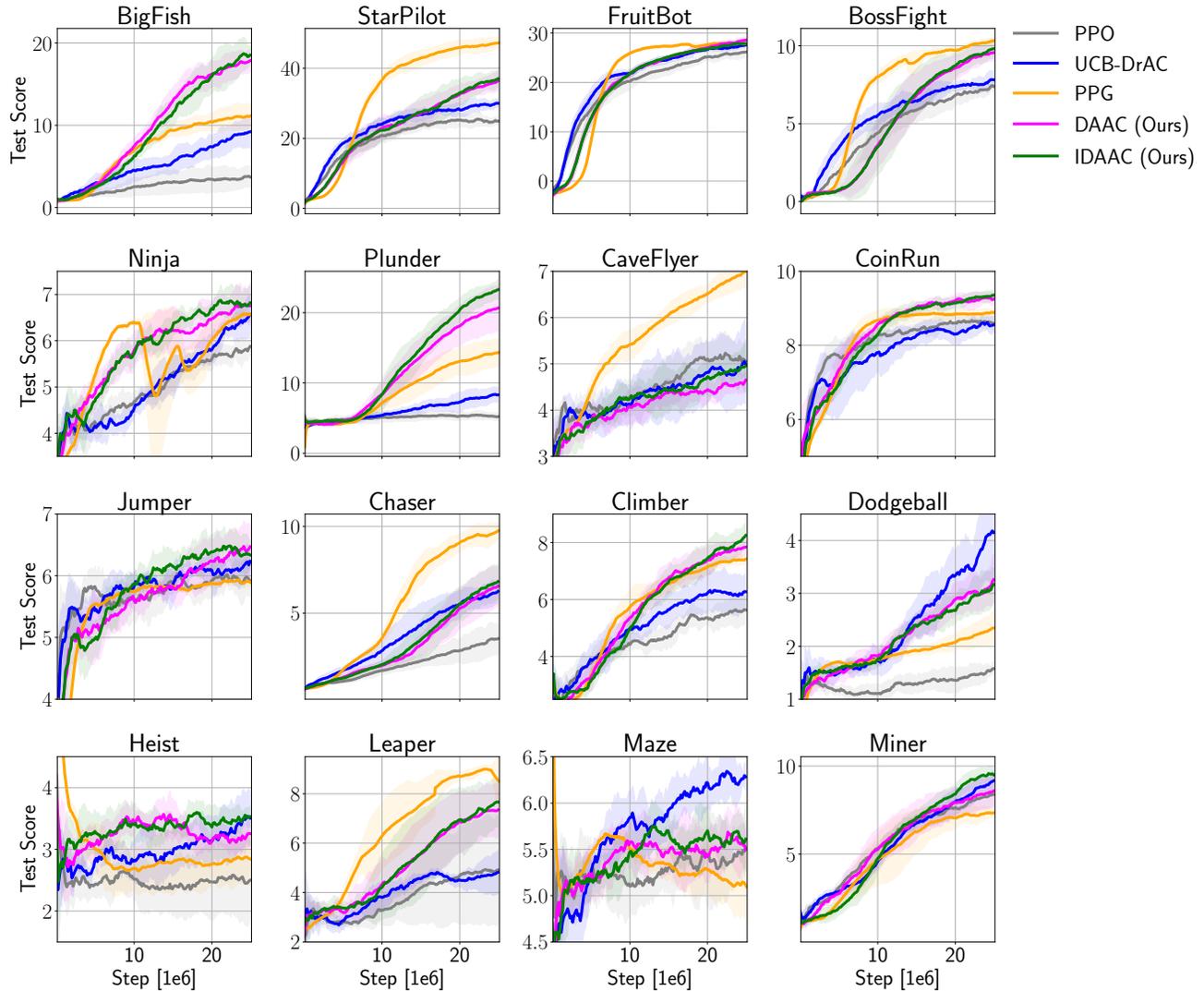


Figure 7. **Test Performance** of IDAAC, DAAC, PPG, UCB-DrAC, and PPO on all Procgen games. IDAAC outperforms the other methods on most games and is significantly better than PPO. The mean and standard deviation are computed over 10 runs with different seeds.

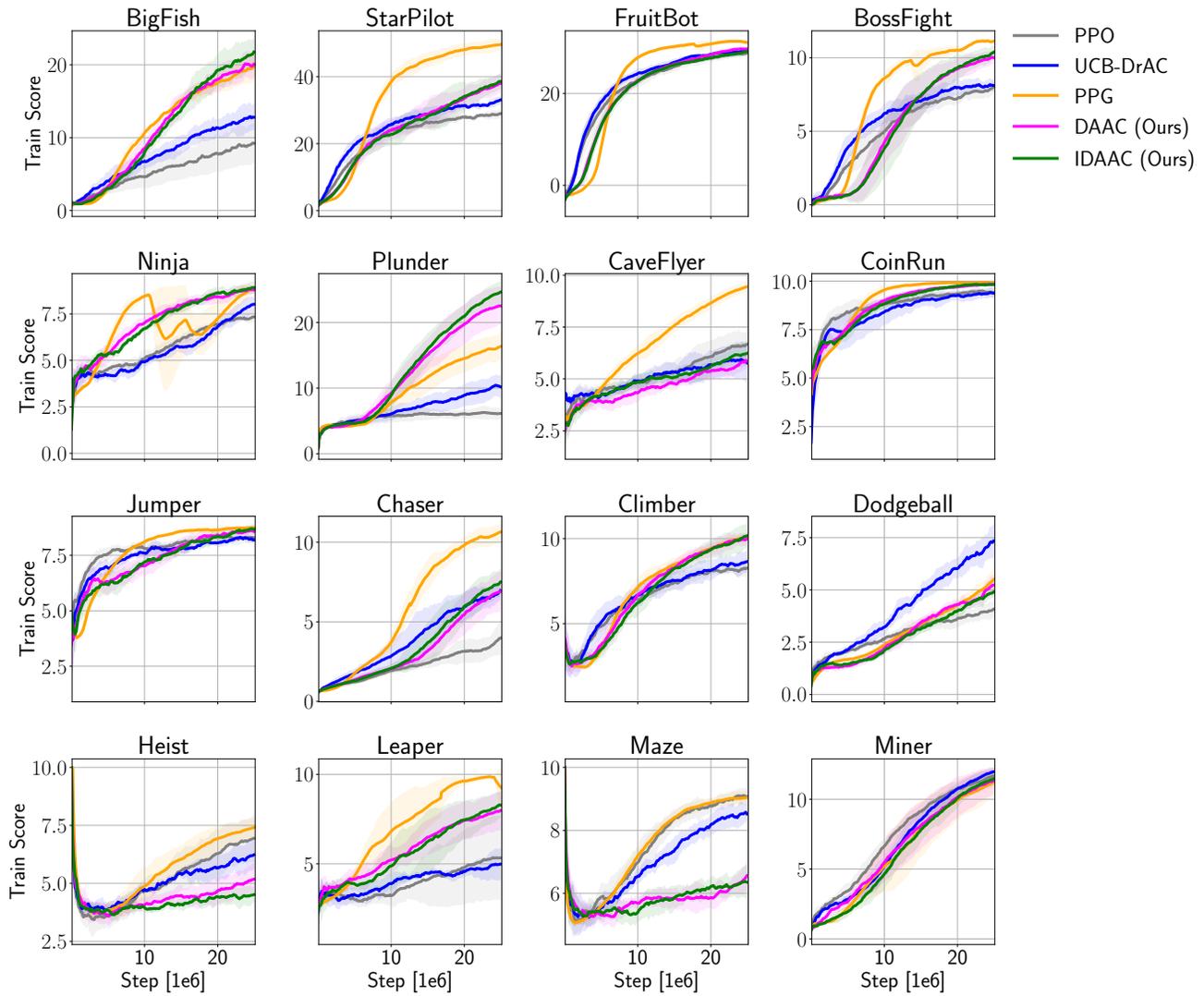


Figure 8. **Train Performance** of IDAAC, DAAC, PPG, UCB-DrAC, and PPO on all Procgen games. IDAAC outperforms the other methods on most games and is significantly better than PPO. The mean and standard deviation are computed over 10 runs with different seeds.

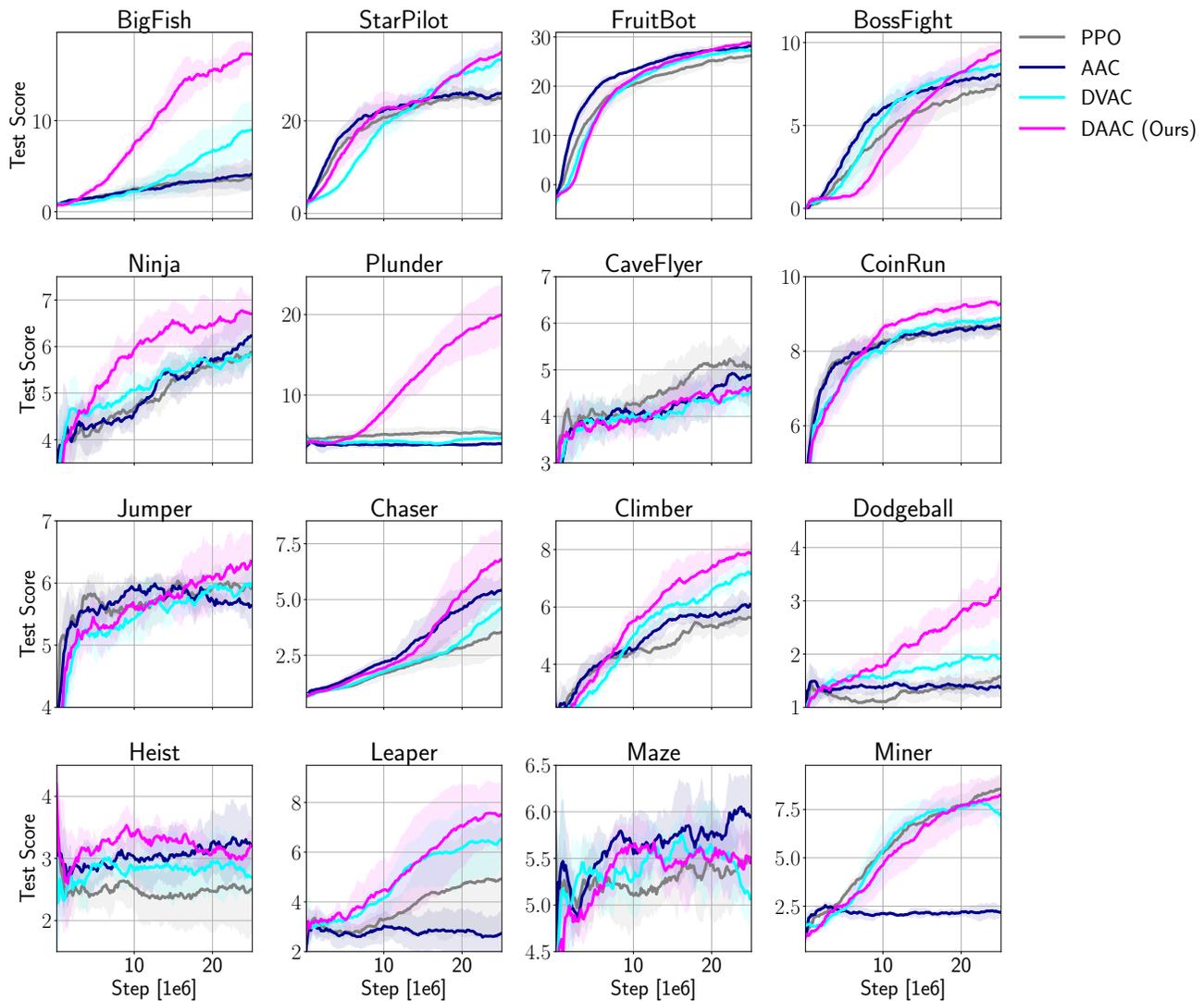


Figure 9. Test Performance of DAAC, DVAC, AAC, and PPO on all Procgen games. DVAC is an ablation of DAAC that replaces the advantage head of the policy network with a value head. AAC is similar to PPO but has an additional advantage head as part of the policy network. DAAC outperforms all these ablations, emphasizing the importance of all of its components. The mean and standard deviation are computed over 5 runs with different seeds.

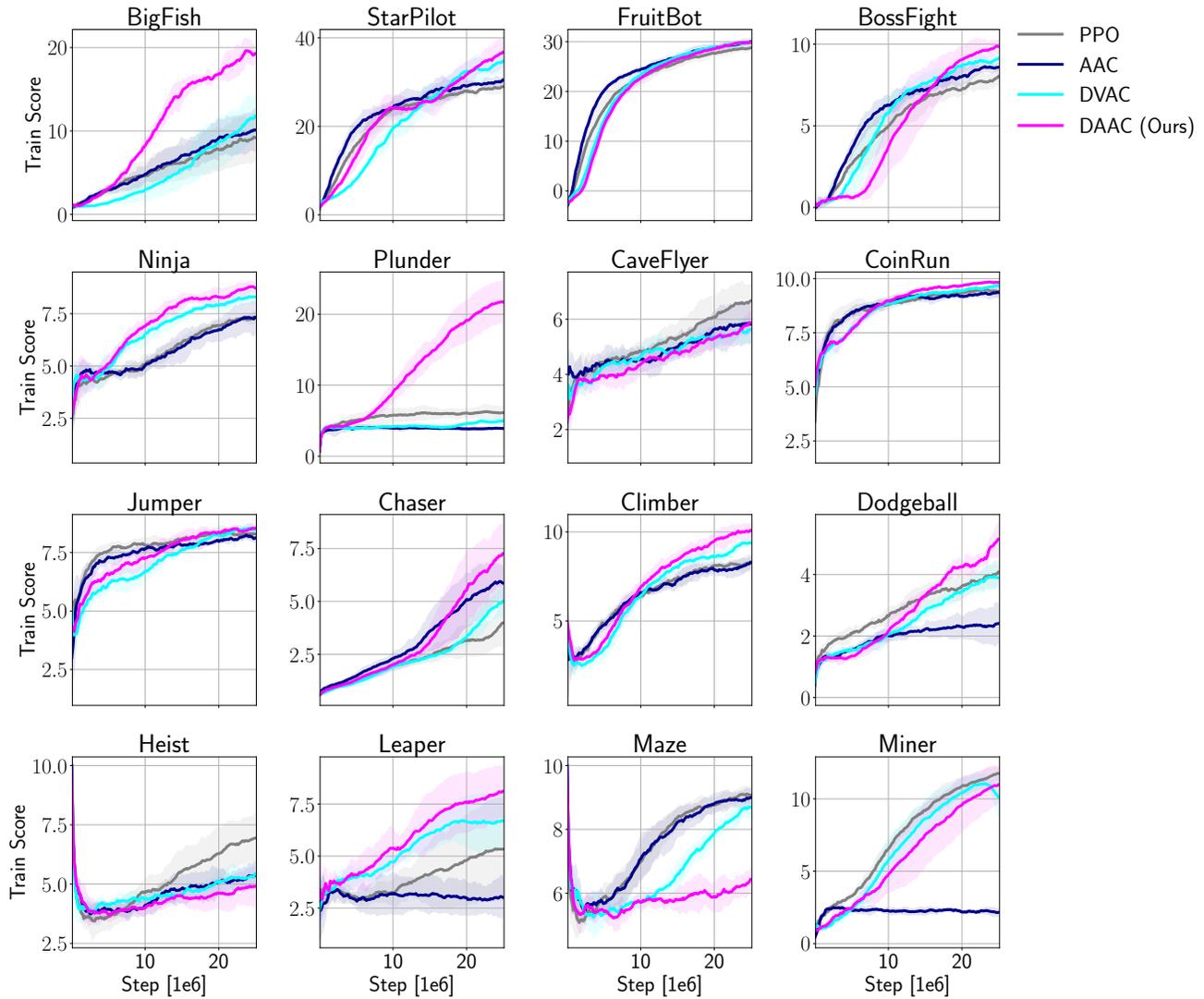


Figure 10. **Train Performance** of DAAC, DVAC, AAC, and PPO on all Procgen games. DVAC is an ablation of DAAC that replaces the advantage head of the policy network with a value head. AAC is similar to PPO but has an additional advantage head as part of the policy network. DAAC outperforms all these ablations, emphasizing the importance of all of its components. The mean and standard deviation are computed over 5 runs with different seeds.

**Decoupling Value and Policy for Generalization in Reinforcement Learning**

Tables 4 and 5 show the final mean and standard deviation of the test and train scores obtained by our methods and the baselines on each of the 16 Procgen games, after 25M training steps. For brevity, we only include the strongest baselines in these tables. See Raileanu et al. (2020) for the scores obtained by Rand-FM and IBAC-SNI.

Table 4. Procgen scores on test levels after training on 25M environment steps. The mean and standard deviation are computed using 10 runs with different seeds.

Game	PPO	Mixreg	PLR	UCB-DrAC	PPG	DAAC (Ours)	IDAAC (Ours)
BigFish	3.7 ± 1.3	7.1 ± 1.6	10.9 ± 2.8	9.2 ± 2.0	11.2 ± 1.4	17.8 ± 1.4	<b>18.5 ± 1.2</b>
StarPilot	24.9 ± 1.0	32.4 ± 1.5	27.9 ± 4.4	30.0 ± 1.3	<b>47.2 ± 1.6</b>	36.4 ± 2.8	37.0 ± 2.3
FruitBot	26.2 ± 1.2	27.3 ± 0.8	28.0 ± 1.4	27.6 ± 0.4	27.8 ± 0.6	<b>28.6 ± 0.6</b>	27.9 ± 0.5
BossFight	7.4 ± 0.4	8.2 ± 0.7	8.9 ± 0.4	7.8 ± 0.6	<b>10.3 ± 0.2</b>	9.6 ± 0.5	9.8 ± 0.6
Ninja	5.9 ± 0.2	6.8 ± 0.5	<b>7.2 ± 0.4</b>	6.6 ± 0.4	6.6 ± 0.1	6.8 ± 0.4	6.8 ± 0.4
Plunder	5.2 ± 0.6	5.9 ± 0.5	8.7 ± 2.2	8.3 ± 1.1	14.3 ± 2.0	20.7 ± 3.3	<b>23.3 ± 1.4</b>
CaveFlyer	5.1 ± 0.4	6.1 ± 0.6	6.3 ± 0.5	5.0 ± 0.8	<b>7.0 ± 0.4</b>	4.6 ± 0.2	5.0 ± 0.6
CoinRun	8.6 ± 0.2	8.6 ± 0.3	8.8 ± 0.5	8.6 ± 0.2	8.9 ± 0.1	9.2 ± 0.2	<b>9.4 ± 0.1</b>
Jumper	5.9 ± 0.2	6.0 ± 0.3	5.8 ± 0.5	6.2 ± 0.3	5.9 ± 0.1	<b>6.5 ± 0.4</b>	6.3 ± 0.2
Chaser	3.5 ± 0.9	5.8 ± 1.1	6.9 ± 1.2	6.3 ± 0.6	<b>9.8 ± 0.5</b>	6.6 ± 1.2	6.8 ± 1.0
Climber	5.6 ± 0.5	6.9 ± 0.7	6.3 ± 0.8	6.3 ± 0.6	2.8 ± 0.4	7.8 ± 0.2	<b>8.3 ± 0.4</b>
Dodgeball	1.6 ± 0.1	1.7 ± 0.4	1.8 ± 0.5	<b>4.2 ± 0.9</b>	2.3 ± 0.3	3.3 ± 0.5	3.2 ± 0.3
Heist	2.5 ± 0.6	2.6 ± 0.4	2.9 ± 0.5	3.5 ± 0.4	2.8 ± 0.4	3.3 ± 0.2	<b>3.5 ± 0.2</b>
Leaper	4.9 ± 2.2	5.3 ± 1.1	6.8 ± 1.2	4.8 ± 0.9	<b>8.5 ± 1.0</b>	7.3 ± 1.1	7.7 ± 1.0
Maze	5.5 ± 0.3	5.2 ± 0.5	5.5 ± 0.8	<b>6.3 ± 0.1</b>	5.1 ± 0.3	5.5 ± 0.2	5.6 ± 0.3
Miner	8.4 ± 0.7	9.4 ± 0.4	9.6 ± 0.6	9.2 ± 0.6	7.4 ± 0.2	8.6 ± 0.9	<b>9.5 ± 0.4</b>

Table 5. Procgen scores on train levels after training on 25M environment steps. The mean and standard deviation are computed using 10 runs with different seeds.

Game	PPO	Mixreg	PLR	UCB-DrAC	PPG	DAAC (Ours)	IDAAC (Ours)
BigFish	9.2 ± 2.7	15.0 ± 1.3	7.8 ± 1.0	12.8 ± 1.8	19.9 ± 1.7	20.1 ± 1.6	<b>21.8 ± 1.8</b>
StarPilot	29.0 ± 1.1	28.7 ± 1.1	2.6 ± 0.3	33.1 ± 1.3	<b>49.6 ± 2.1</b>	38.0 ± 2.6	38.6 ± 2.2
FruitBot	28.8 ± 0.6	29.9 ± 0.5	15.9 ± 1.3	29.3 ± 0.5	<b>31.1 ± 0.5</b>	29.7 ± 0.4	29.1 ± 0.7
BossFight	8.0 ± 0.4	7.9 ± 0.8	8.7 ± 0.7	8.1 ± 0.4	<b>11.1 ± 0.1</b>	10.0 ± 0.4	10.4 ± 0.4
Ninja	7.3 ± 0.3	8.2 ± 0.4	5.4 ± 0.5	8.0 ± 0.4	8.9 ± 0.2	8.8 ± 0.2	<b>8.9 ± 0.3</b>
Plunder	6.1 ± 0.8	6.2 ± 0.3	4.1 ± 1.3	10.2 ± 1.76	16.4 ± 1.9	22.5 ± 2.8	<b>24.6 ± 1.6</b>
CaveFlyer	6.7 ± 0.6	6.2 ± 0.7	6.4 ± 0.1	5.8 ± 0.9	<b>9.5 ± 0.2</b>	5.8 ± 0.4	6.2 ± 0.6
CoinRun	9.4 ± 0.3	9.5 ± 0.2	5.4 ± 0.4	9.4 ± 0.2	<b>9.9 ± 0.0</b>	9.8 ± 0.0	9.8 ± 0.1
Jumper	8.3 ± 0.2	8.5 ± 0.4	3.6 ± 0.5	8.2 ± 0.1	8.7 ± 0.1	8.6 ± 0.3	<b>8.7 ± 0.2</b>
Chaser	4.1 ± 0.3	3.4 ± 0.9	6.3 ± 0.7	7.0 ± 0.6	<b>10.7 ± 0.4</b>	6.9 ± 1.2	7.5 ± 0.8
Climber	6.9 ± 1.0	7.5 ± 0.8	6.2 ± 0.8	8.6 ± 0.6	10.2 ± 0.2	10.0 ± 0.3	<b>10.2 ± 0.7</b>
Dodgeball	5.3 ± 2.3	9.1 ± 0.5	2.0 ± 1.1	<b>7.3 ± 0.8</b>	5.5 ± 0.5	5.2 ± 0.4	4.9 ± 0.3
Heist	7.1 ± 0.5	4.4 ± 0.3	1.2 ± 0.4	6.2 ± 0.6	<b>7.4 ± 0.4</b>	5.2 ± 0.7	4.5 ± 0.3
Leaper	5.5 ± 0.4	3.2 ± 1.2	6.4 ± 0.4	5.0 ± 0.9	<b>9.3 ± 1.1</b>	8.0 ± 1.1	8.3 ± 0.7
Maze	<b>9.1 ± 0.2</b>	8.7 ± 0.7	4.1 ± 0.5	8.5 ± 0.3	9.0 ± 0.2	6.6 ± 0.4	6.4 ± 0.5
Miner	11.7 ± 0.5	8.9 ± 0.9	9.7 ± 0.4	<b>12.0 ± 0.3</b>	11.3 ± 1.0	11.3 ± 0.9	11.5 ± 0.5

## E. DeepMind Control Suite Experiments

For our DMC experiments, we followed the protocol proposed in Zhang et al. (2020b) to modify the tasks so that they contain natural and synthetic distractors in the background. For each DMC task, we split the generated environments (each with a different background video) into training (80%) and testing (20%). The results shown here correspond to the average return on the test environments, over the course of training. Note that this setting is slightly different from the one used in Raileanu et al. (2020) which shows results on all the generated environments, just like Zhang et al. (2020b). As Figure 11 shows, our methods outperform the baselines on these continuous control tasks.

In line with standard practice for this benchmark, we use 8 action repeats for Cartpole Swingup and 4 for Cartpole Balance and Ball In Cup. We also use 3 stacked frames as observations. To find the best hyperparameters, we ran a grid search over the learning rate in  $[0.0001, 0.0003, 0.0007, 0.001]$ , the number of minibatches in  $[32, 8, 16, 64]$ , the entropy coefficient in  $[0.0, 0.01, 0.001, 0.0001]$ , and the number of PPO epochs per update in  $[3, 5, 10, 20]$ . We found 10 ppo epochs, 0.0 entropy coefficient, 0.0003 learning rate, and 32 minibatches to work best across these environments. We use  $\gamma = 0.99$ ,  $\lambda = 0.95$  for the generalized advantage estimates, 2048 steps, 1 process, value loss coefficient 0.5, and linear rate decay over 1 million environment steps. Following this grid search, we used the best values found for all the methods. Any other hyperparameters not mentioned here were set to the same values as the ones used for Procgen as described above. For UCB-DrAC, we used the best hyperparameters found by the corresponding authors. For PPG, we ran the same hyperparameter search as the one performed in the original paper for Procgen and found  $N_\pi = 32$ ,  $E_\pi = 1$ ,  $E_V = 1$ ,  $E_{aux} = 6$ , and  $\beta_{clone} = 1$  to be the best. Similarly, for DAAC and IDAAC, we ran the same hyperparameter search as for Procgen and found that  $E_V = 9$ ,  $N_\pi = 32$ ,  $\alpha_a = 0.1$ , and  $\alpha_i = 0.1$  worked best across all environments.

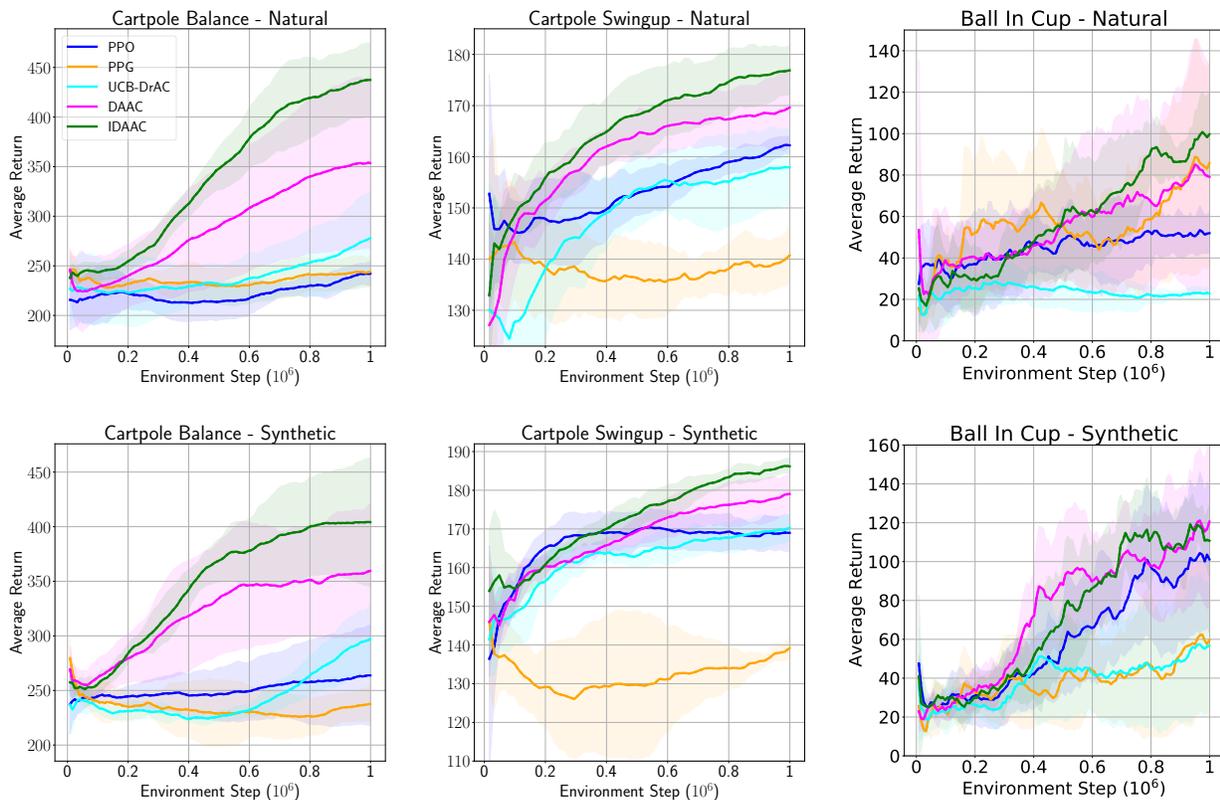


Figure 11. Average return on three DMC tasks, Cartpole Balance (left), Cartpole Swingup (center), and Ball In Cup (right), with natural (top) and synthetic (bottom) video backgrounds. The mean and standard deviation are computed over 10 runs with different seeds. DAAC and IDAAC outperform PPO, PPG, and UCB-DrAC.

## F. Value Loss and Generalization

In this section, we look at the relationship between the value loss, test score, and number of training levels for all Progen games (see Figures 12, 14, and 13). As discussed in the paper, the value loss is *positively* correlated with the test score and number of training levels. This result goes against our intuition from training RL algorithms on single environments where in general, the value loss is *inversely* correlated with the agent’s performance and sample efficiency. This observation further supports our claim that having a more accurate value function can lead to representations that overfit to the training environments. When using a shared network for the policy and value, this results in policies that do not generalize well to new environments. By decoupling the representations of the policy and value function, our methods DAAC and IDAAC can achieve the best of both worlds by (i) learning accurate value functions, while also (ii) learning representations and policies that better generalize to unseen environments.

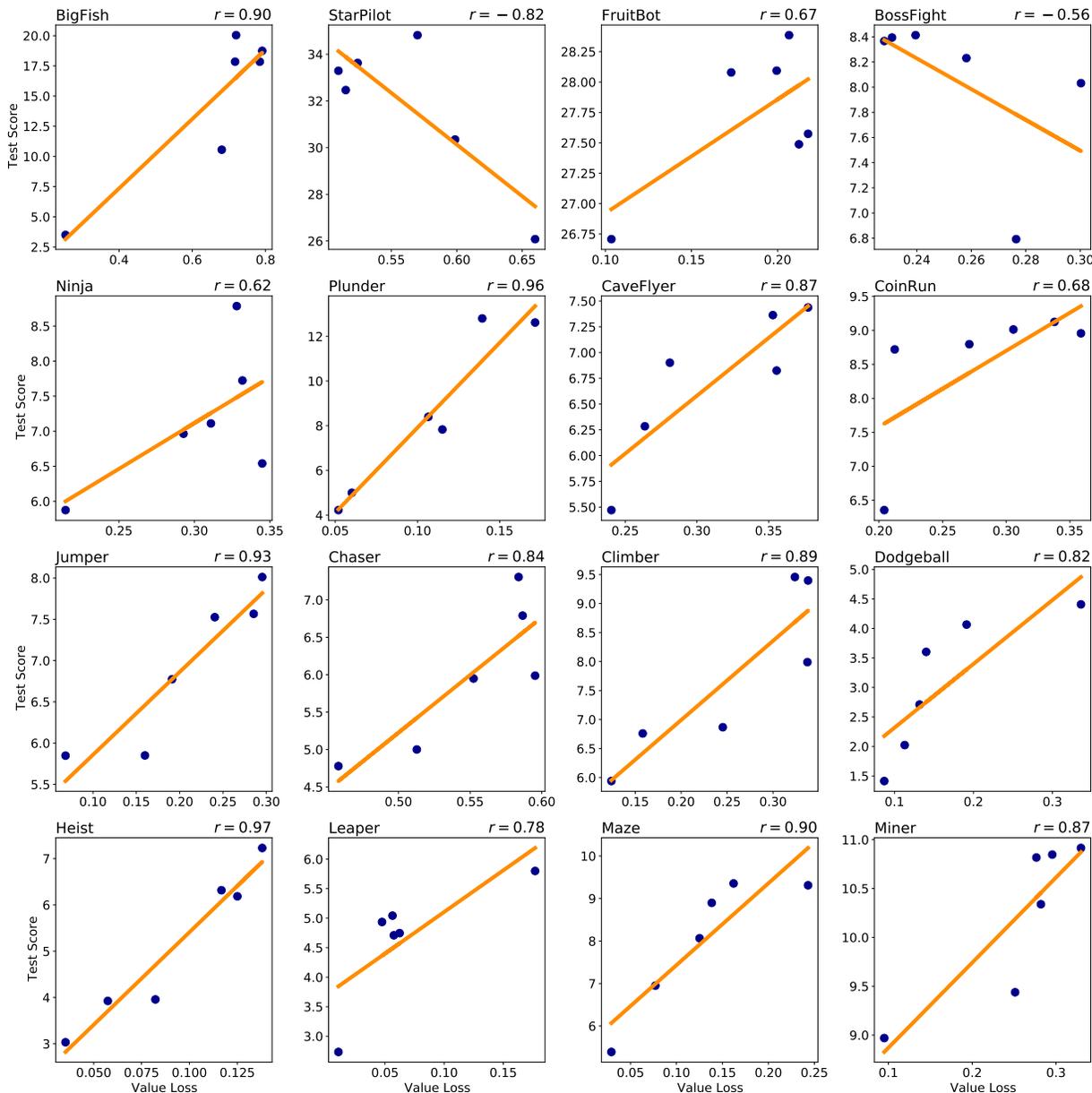


Figure 12. Correlation of the value loss and test score for PPO agents trained on varying numbers of levels: 200, 500, 1000, 2000, 5000, and 10000. Surprisingly, the value loss is positively correlated with the test score for most games, suggesting that *models with larger value loss generalize better*.

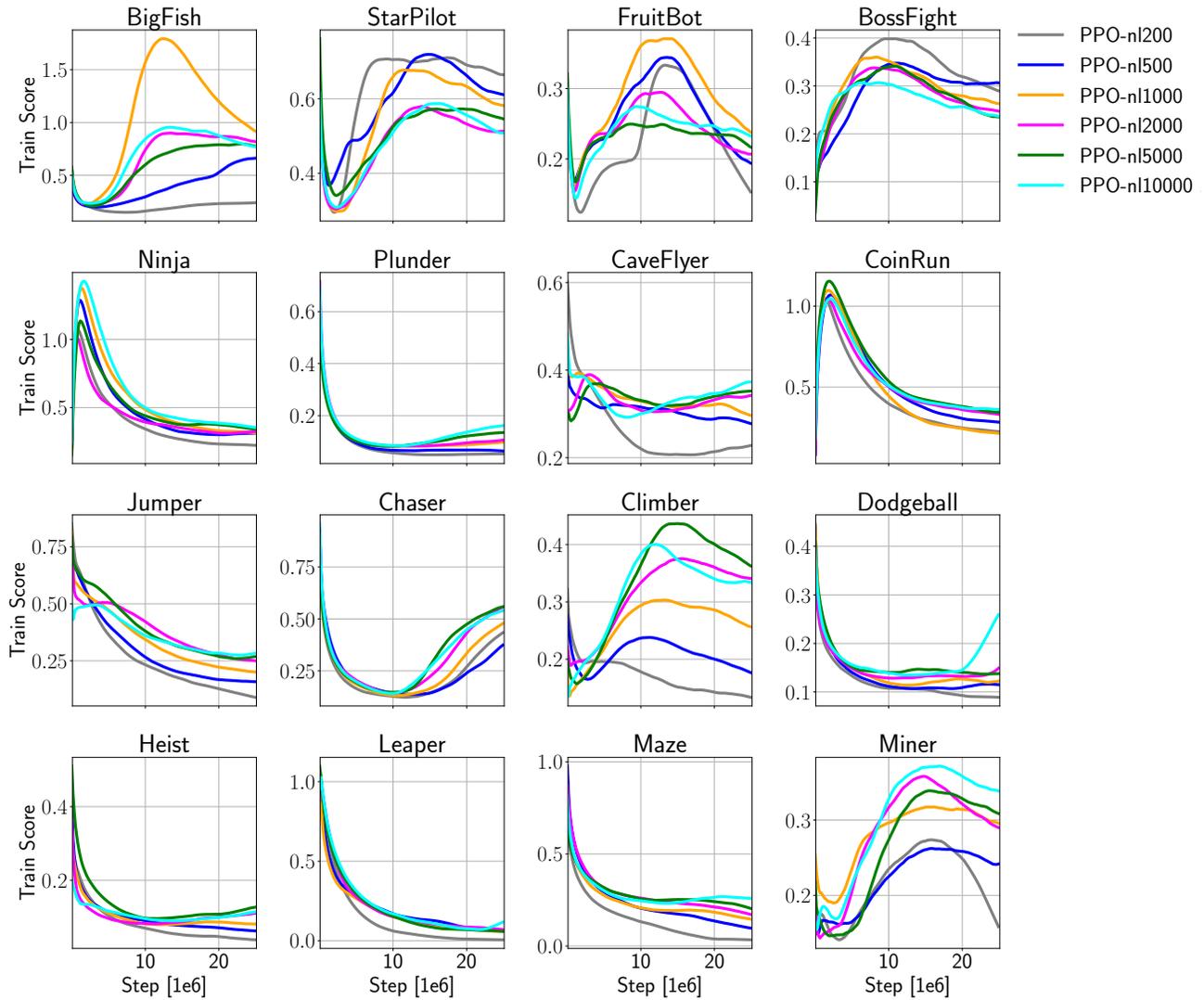


Figure 13. Test score for PPO agents trained on varying numbers of levels: 200, 500, 1000, 2000, 5000, and 10000. For most games, the test score increases with the number of training levels, suggesting that models trained on more levels generalize better to unseen levels, as expected.

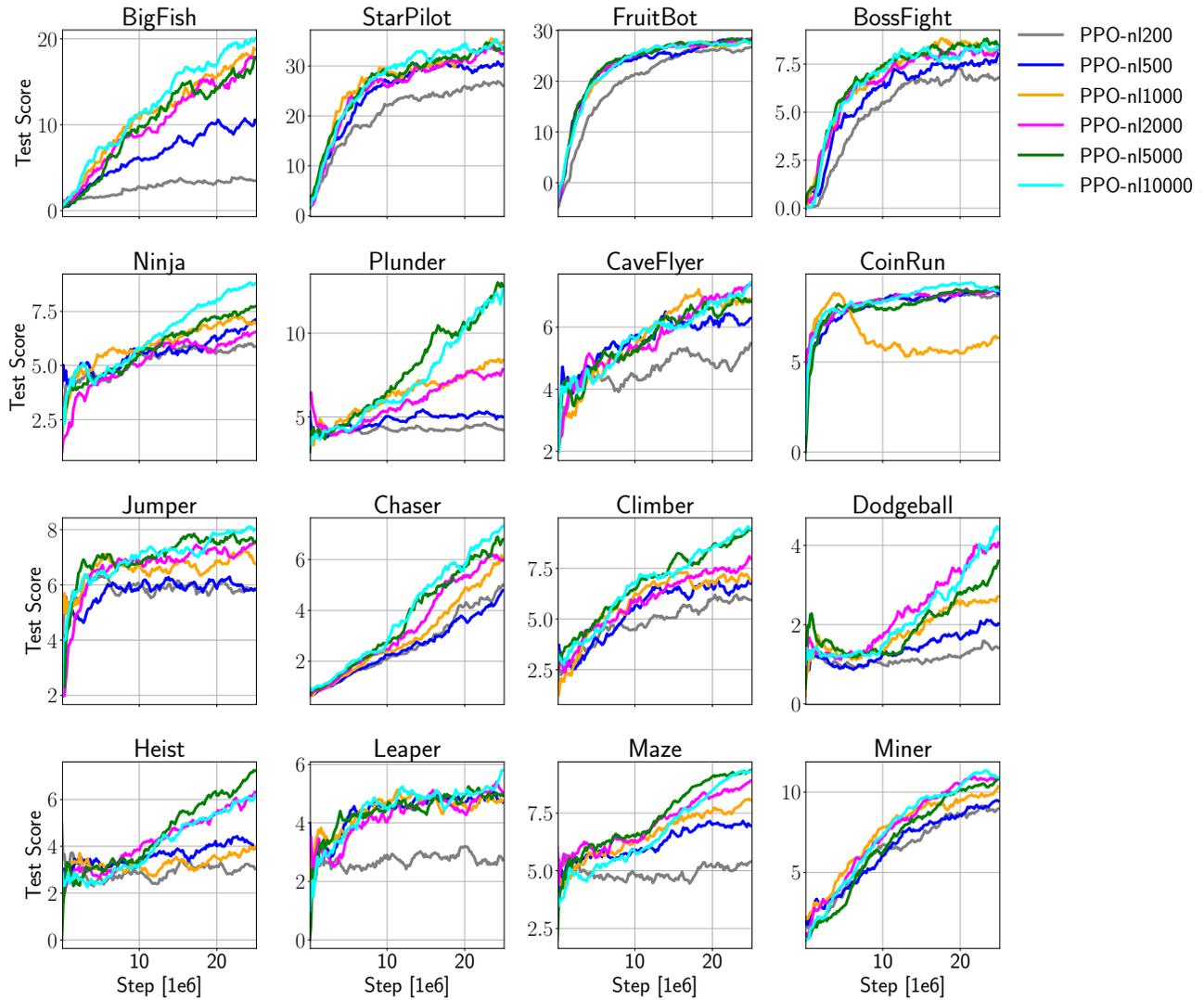


Figure 14. Value loss for PPO agents trained on varying numbers of levels: 200, 500, 1000, 2000, 5000, and 10000. For most games, the value loss increases with the number of training levels used, suggesting that models trained on more levels (and thus with better generalization) have higher value loss.

### G. Advantage vs. Value During an Episode

Figure 15 shows two different Ninja levels and their corresponding true and predicted values and advantages (by a PPO agent trained on 200 levels). The true values and advantages are computed assuming an optimal policy. For illustration purposes, we show the advantage for the noop action, but similar conclusions apply to the other actions. As the figure shows, the true value function is different for the two levels, while the advantage is (approximately) the same. This holds true for the agent’s estimates as well, suggesting that using the value loss to update the policy parameters can lead to more overfitting than using the advantage loss (since it exhibits less dependence on the idiosyncrasies of a level such as its length or difficulty).

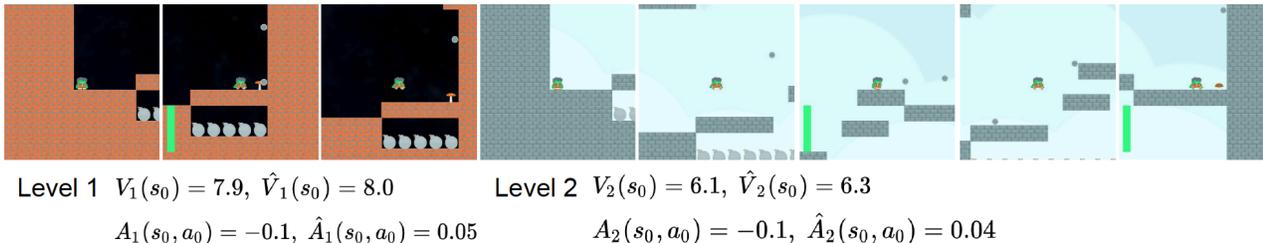


Figure 15. **Policy-Value Asymmetry.** Two Ninja levels with initial observations that are *semantically identical but visually different*. Level 1 (first three frames from the left with black background) is much shorter than Level 2 (last six frames with blue background). Both the true and the estimated values (by a PPO agent trained on 200 levels) of the initial observation are higher for Level 1 than for Level 2 *i.e.*  $V_1(s_0) > V_2(s_0)$  and  $\hat{V}_1(s_0) > \hat{V}_2(s_0)$ . Thus, to accurately predict the value function, the representations must capture level-specific features (such as the backgrounds), which are irrelevant for finding the optimal policy. Consequently, using a common representation for both the policy and value function can lead to overfitting to spurious correlations and poor generalization to unseen levels. In contrast to the value, the true advantage of the initial states and noop action has the same values for the two levels, and the advantages predicted by the agent also have very similar values.

We also analyze the timestep-dependence of the predictions (*i.e.* value or advantage) made by various models over the course of an episode. Figures 16, 17, 18, and 19 show examples from CoinRun, Ninja, Climber, and Jumper, respectively. While both PPO and PPG (trained on 200 levels) learn a value function which is increasing almost linearly with the episode step, DAAC (also trained on 200 levels) learns an advantage function which does not have a clear dependence on the episode step. Thus, DAAC is less prone to overfitting than PPO or PPG. Similarly, a PPO model trained on 10k levels with good generalization performance does not display a linear trend between value and episode step. This suggests there is a trade-off between value accuracy and generalization performance. However, by decoupling the policy and value representations, DAAC is able to achieve both accurate value predictions and good generalization abilities.

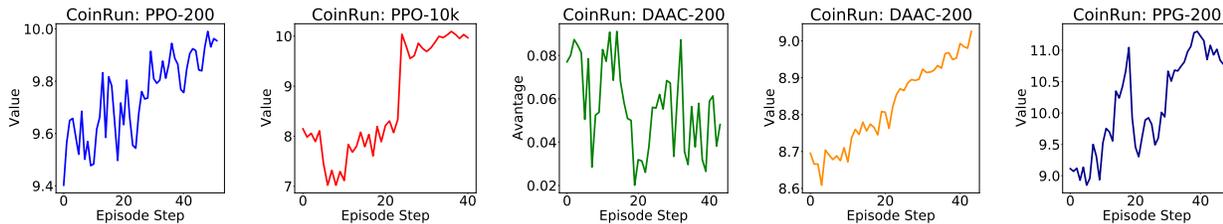


Figure 16. **Examples illustrating the timestep-dependence of the value function for a single CoinRun level.** The dark and light blue curves show the value as a function of episode step for PPG and PPO, respectively, each trained on 200 levels. Note the near linear relationship, indicating overfitting to the training levels. By contrast, a PPO model (red) trained on 10k levels (thus exhibiting far less overfitting) does not show this relationship. Our DAAC model trained on 200 levels (green) also lacks this adverse dependence in the advantage prediction which is used for training the policy network, thus is able to generalize better than the PPO model trained on the same amount of data (see Fig. 3). Nevertheless, DAAC’s value estimate still have a linear trend (orange) but, in contrast to PPO and PPG, this does not negatively affect the policy since we use separate networks for learning the policy and value.

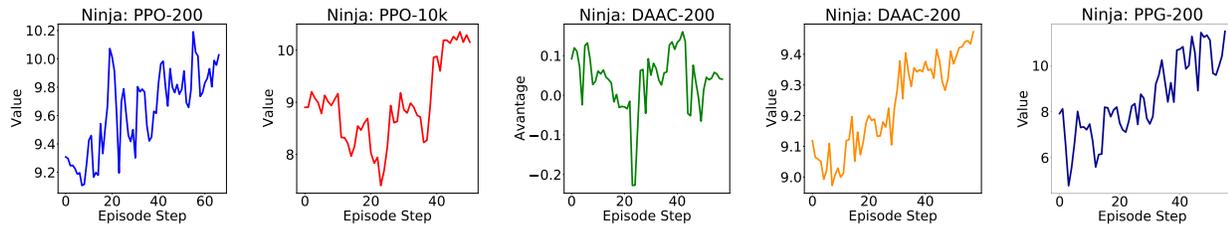


Figure 17. Examples illustrating the timestep-dependence of the value function for a single Ninja level. The dark and light blue curves show the value as a function of episode step for PPG and PPO, respectively, each trained on 200 levels. Note the near linear relationship, indicating overfitting to the training levels. By contrast, a PPO model (red) trained on 10k levels (thus exhibiting far less overfitting) does not show this relationship. Our DAAC model trained on 200 levels (green) also lacks this adverse dependence in the advantage prediction which is used for training the policy network, thus is able to generalize better than the PPO model trained on the same amount of data (see Fig. 3). Nevertheless, DAAC’s value estimate still have a linear trend (orange) but, in contrast to PPO and PPG, this does not negatively affect the policy since we use separate networks for learning the policy and value.

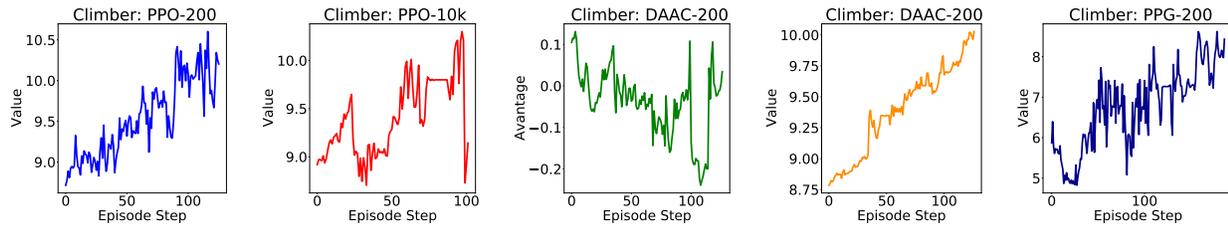


Figure 18. Examples illustrating the timestep-dependence of the value function for a single Climber level. The dark and light blue curves show the value as a function of episode step for PPG and PPO, respectively, each trained on 200 levels. Note the near linear relationship, indicating overfitting to the training levels. By contrast, a PPO model (red) trained on 10k levels (thus exhibiting far less overfitting) does not show this relationship. Our DAAC model trained on 200 levels (green) also lacks this adverse dependence in the advantage prediction which is used for training the policy network, thus is able to generalize better than the PPO model trained on the same amount of data (see Fig. 3). Nevertheless, DAAC’s value estimate still have a linear trend (orange) but, in contrast to PPO and PPG, this does not negatively affect the policy since we use separate networks for learning the policy and value.

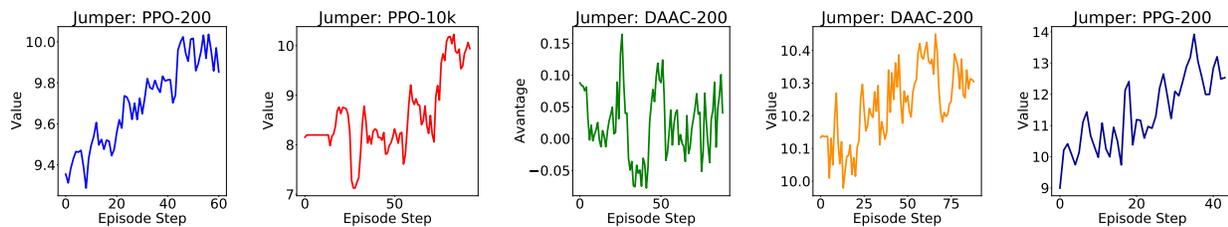


Figure 19. Examples illustrating the timestep-dependence of the value function for a single Jumper level. The dark and light blue curves show the value as a function of episode step for PPG and PPO, respectively, each trained on 200 levels. Note the near linear relationship, indicating overfitting to the training levels. By contrast, a PPO model (red) trained on 10k levels (thus exhibiting far less overfitting) does not show this relationship. Our DAAC model trained on 200 levels (green) also lacks this adverse dependence in the advantage prediction which is used for training the policy network, thus is able to generalize better than the PPO model trained on the same amount of data (see Fig. 3). Nevertheless, DAAC’s value estimate still have a linear trend (orange) but, in contrast to PPO and PPG, this does not negatively affect the policy since we use separate networks for learning the policy and value.

## H. Value Variance

In this section, we look at the variance in the predicted values for the initial observation, across all training levels. In partially-observed procedurally generated environments, there should be no way of telling how difficult or long a level is (and thus how much reward is expected) from the initial observation alone since the end of the level cannot be seen. Thus, we would expect a model with strong generalization to predict similar values for the initial observation irrespective of the environment instance. If this is not the case and the model uses a common representation for the policy and value function, the policy is likely to overfit to the training environments. As Figure 20 shows, the variance decreases with the number of levels used for training. This is consistent with our observation that models with better generalization predict more similar values for observations that are semantically different such as the initial observation. In contrast, models trained on a low number of levels, memorize the value of the initial observation for each level, leading to poor generalization in new environments. Note that we chose to illustrate this phenomenon on three of the Procgen games (*i.e.* Climber, Jumper, and Ninja) where it is more apparent due to their partial-observability and substantial level diversity (in terms of length).

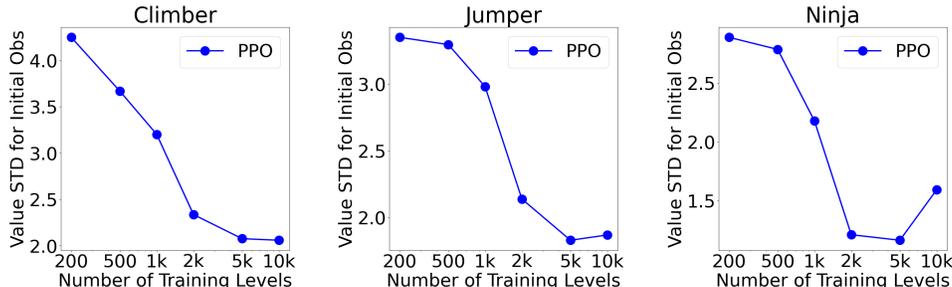


Figure 20. Standard deviation of the predicted values for the initial observations from 200 levels, as a function of the number of training levels. From left to right: Climber, Jumper, and Ninja. The 200 levels used to compute the standard deviation were part of the training set for all the agents. Note that, in general, the variance of the value decreases with the number of training levels. This is consistent with our claim that when sharing parameters for the policy and value, models which generalize better predict close values for observations that are semantically similar (*e.g.* the initial observation) since they learn representations which are less prone to overfitting to level-specific features.

## I. Robustness to Spurious Correlations

In this section, we investigate how robust the learned features, policies, and predicted values or advantages, are to spurious correlations. To answer this question, we measure how much the features, policies, and predictions vary when the background of an observation changes. Note that the change in background does not change the underlying state of the environment but only its visual aspect. Hence, a change in background should not modify the agent’s policy or learned representation. For this experiment, we collect a buffer of 1000 observations from 20 training levels, using a PPO agent trained on 200 levels. Then, we create 10 extra versions of each observation by changing the background. We then measure the L1-norm and L2-norm between the learned representation (*i.e.* final vector before the policy’s softmax layer) of the original observation and each of its other versions. We also compute the difference in predicted outputs (*i.e.* values for PPO and PPG or advantages for DAAC and IDAAC) and the Jensen-Shannon Divergence (JSD) between the policies. For all these metrics, we first take the mean for all 10 backgrounds to obtain a single point for each observation, and then we report the mean and standard deviation across all 1000 observations, resulting in an average statistic of how much these metrics change as a result of varying the background.

Figures 21, 22, and 23 show the results for Ninja, Jumper, and Climber, respectively, comparing IDAAC, DAAC, PPG, as well as PPO trained on 200 and 10k levels. In particular, the results show that both our methods learn representations which are more robust to changes in the background than PPO and PPG (assuming all methods are trained on the same number of levels *i.e.* 200). Overall, the differences due to background changes in the auxiliary outputs of the policy networks for DAAC and IDAAC (*i.e.* the predicted advantages) are smaller than those of PPO and PPG (*i.e.* the predicted values). These results indicate that DAAC and IDAAC are more robust than PPO and PPG to visual features which are irrelevant for control.

We do not observe a significant difference across the JSDs of the different methods. However, in the case of Procgen, the

JSD isn't a perfect measure of the semantic difference between two policies because some of the actions have the same effect on the environment (*e.g.* in Ninja, there are two actions that move the agent to the right) and thus are interchangeable. Two policies could have a large JSD while being semantically similar, thus rendering the policy robustness analysis inconclusive.

In some cases, PPO-10k exhibits better robustness to different backgrounds than DAAC and IDAAC by exhibiting a lower feature norm and value or advantage difference. However, PPO-10k is a PPO model trained on 10000 levels of a game, while DAAC and IDAAC are trained only on 200 levels. For most Procgen games, training on 10k levels is enough to generalize to the test distribution so PPO-10k is expected to generalize better than methods trained on 200 levels. In this paper, we are interested in generalizing to unseen levels from a small number of training levels, so PPO-10k is used as an upper-bound rather than a baseline since a direct comparison wouldn't be fair. Hence, it is not surprising that some of these robustness metrics are better for PPO-10k than DAAC and IDAAC.

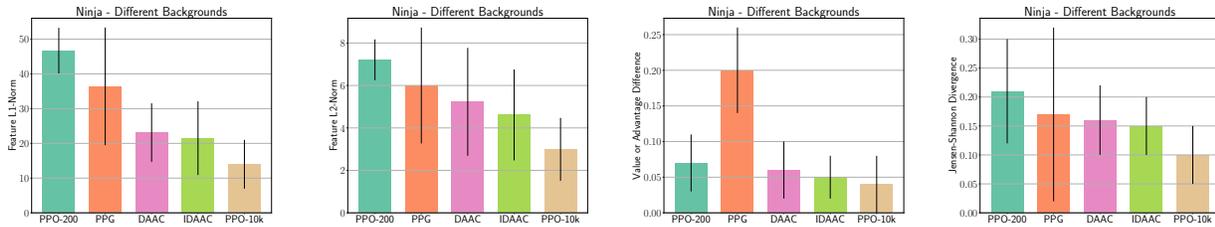


Figure 21. Variations in the learned features, policies, and values or advantages when changing the background in Ninja. From left to right we report the L1 and L2-norm for the features, the value or advantage difference, and the Jensen-Shannon Divergence for the policy. We compare PPO trained on 200 and 10k levels with PPG, DAAC, and IDAAC. Our models are more robust to changes in the background (which does not affect the state). The means and standard deviations were computed over 10 different backgrounds.

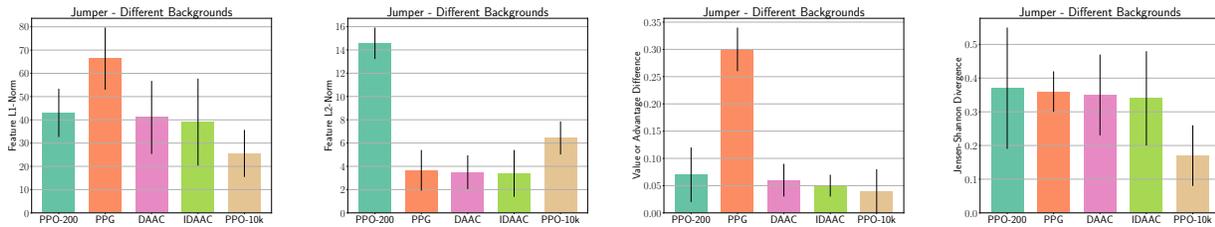


Figure 22. Variations in the learned features, policies, and values or advantages when changing the background in Jumper. From left to right we report the L1 and L2-norm for the features, the value or advantage difference, and the Jensen-Shannon Divergence for the policy. We compare PPO trained on 200 and 10k levels with PPG, DAAC, and IDAAC. Our models are more robust to changes in the background (which does not affect the state). The means and standard deviations were computed over 10 different backgrounds.

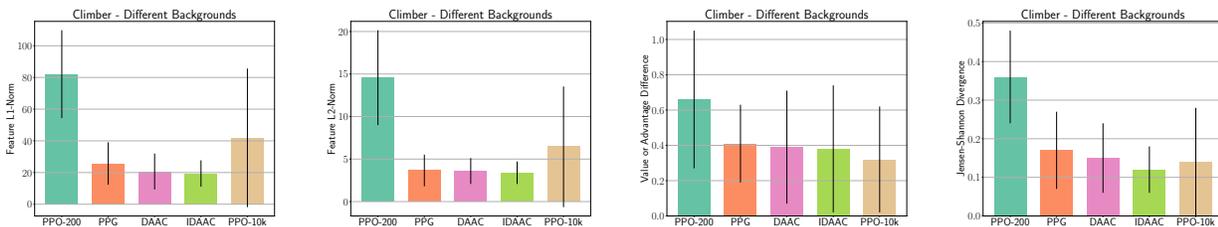


Figure 23. Variations in the learned features, policies, and values or advantages when changing the background in Climber. From left to right we report the L1 and L2-norm for the features, the value or advantage difference, and the Jensen-Shannon Divergence for the policy. We compare PPO trained on 200 and 10k levels with PPG, DAAC, and IDAAC. Our models are more robust to changes in the background (which does not affect the state). The means and standard deviations were computed over 10 different backgrounds.