

# Improving Sample Efficiency in Model-Free Reinforcement Learning from Images

Denis Yarats<sup>1,2</sup> Amy Zhang<sup>3,4,2</sup> Ilya Kostrikov<sup>1</sup> Brandon Amos<sup>2</sup> Joelle Pineau<sup>3,4,2</sup> Rob Fergus<sup>1,2</sup>

## Abstract

Training an agent to solve control tasks directly from high-dimensional images with model-free reinforcement learning (RL) has proven difficult. A promising approach is to learn a latent representation together with the control policy. However, fitting a high-capacity encoder using a scarce reward signal is sample inefficient and leads to poor performance. Prior work has shown that auxiliary losses, such as image reconstruction, can aid efficient representation learning. However, incorporating reconstruction loss into an off-policy learning algorithm often leads to training instability. We explore the underlying reasons and identify variational autoencoders, used by previous investigations, as the cause of the divergence. Following these findings, we propose effective techniques to improve training stability. This results in a simple approach capable of matching state-of-the-art model-free and model-based algorithms on MuJoCo control tasks. Furthermore, our approach demonstrates robustness to observational noise, surpassing existing approaches in this setting. Code, results, and videos are anonymously available at <https://sites.google.com/view/sac-ae/home>.

## 1. Introduction

Cameras are a convenient and inexpensive way to acquire state information, especially in complex, unstructured environments, where effective control requires access to the proprioceptive state of the underlying dynamics. Thus, having effective RL approaches that can utilize pixels as input would potentially enable solutions for a wide range of real world applications, for example robotics.

The challenge is to efficiently learn a mapping from pix-

<sup>1</sup>New York University <sup>2</sup>Facebook AI Research <sup>3</sup>McGill University <sup>4</sup>MILA. Correspondence to: Denis Yarats <denisyarats@cs.nyu.edu>.

els to an appropriate representation for control using only a sparse reward signal. Although deep convolutional encoders can learn good representations (upon which a policy can be trained), they require large amounts of training data. As existing reinforcement learning approaches already have poor sample complexity, this makes direct use of pixel-based inputs prohibitively slow. For example, model-free methods on Atari (Bellemare et al., 2013) and DeepMind Control (DMC) (Tassa et al., 2018) take tens of millions of steps (Mnih et al., 2013; Barth-Maron et al., 2018), which is impractical in many applications, especially robotics.

Some natural solutions to improve sample efficiency are i) to use off-policy methods and ii) add an auxiliary task with an unsupervised objective. Off-policy methods enable more efficient sample re-use, while the simplest auxiliary task is an autoencoder with a pixel reconstruction objective. Prior work has attempted to learn state representations from pixels with autoencoders, utilizing a two-step training procedure, where the representation is first trained via the autoencoder, and then either with a policy learned on top of the fixed representation (Lange & Riedmiller, 2010; Munk et al., 2016; Higgins et al., 2017b; Zhang et al., 2018a; Nair et al., 2018; Dwibedi et al., 2018), or with planning (Matten et al., 2012; Finn et al., 2015). This allows for additional stability in optimization by circumventing dueling training objectives but leads to suboptimal policies. Other work utilizes continual model-free learning with an auxiliary reconstruction signal in an on-policy manner (Jaderberg et al., 2017; Shelhamer et al., 2016). However, these methods do not report of learning representations and a policy jointly in the off-policy setting, or note that it performs poorly (Shelhamer et al., 2016).

We revisit the concept of adding an autoencoder to model-free RL approaches, with a focus on *off-policy* algorithms. We perform a sequence of careful experiments to understand why previous approaches did not work well. We confirm that a pixel reconstruction loss is vital for learning a good representation, specifically when trained jointly, but requires careful design choices to succeed. Based on these findings, we recommend a simple and effective autoencoder-based *off-policy* method that can be trained

*end-to-end*. We believe this to be the first *model-free off-policy* approach to train the latent state representation and policy *jointly* and match performance with state-of-the-art model-based methods<sup>1</sup> (Hafner et al., 2018; Lee et al., 2019) on many challenging control tasks. In addition, we demonstrate robustness to observational noise and outperform prior methods in this more practical setup.

This paper makes three main contributions: (i) a methodical study of the issues involved with combining autoencoders with model-free RL in the off-policy setting that advises a successful variant we call **SAC+AE**; (ii) a demonstration of the robustness of our model-free approach over model-based methods on tasks with noisy observations; and (iii) an open-source PyTorch implementation of our simple and effective algorithm for researchers and practitioners to build upon.

## 2. Related Work

Efficient learning from high-dimensional pixel observations has been a problem of paramount importance for model-free RL. While some impressive progress has been made applying model-free RL to domains with simple dynamics and discrete action spaces (Mnih et al., 2013), attempts to scale these approaches to complex continuous control environments have largely been unsuccessful, both in simulation and the real world. A glaring issue is that the RL signal is much sparser than in supervised learning, which leads to sample inefficiency, and higher dimensional observation spaces such as pixels worsens this problem.

One approach to alleviate this problem is by training with auxiliary losses. Early work (Lange & Riedmiller, 2010) explores using deep autoencoders to learn feature spaces in visual reinforcement learning, crucially Lange & Riedmiller (2010) propose to recompute features for all collected experiences after each update of the autoencoder, rendering this approach impractical to scale to more complicated domains. Moreover, this method has been only demonstrated on toy problems. Alternatively, Finn et al. (2015) apply deep autoencoder pretraining to real world robots that does not require iterative re-training, improving upon computational complexity of earlier methods. However, in this work the linear policy is trained separately from the autoencoder, which we find to not perform as well as end-to-end methods.

Shelhamer et al. (2016) employ auxiliary losses to enhance performance of A3C (Mnih et al., 2016) on Atari. They recommend a multi-task setting and learning dynamics and reward to find a good representation, which relies on the

<sup>1</sup>We define model-based methods as those that train a dynamics model. By this definition, SLAC (Lee et al., 2019) is a model-based method.

assumption that the dynamics in the task are easy to learn and useful for learning a good policy. To prevent instabilities in learning, Shelhamer et al. (2016) pre-train the agent on randomly collected transitions and then perform joint optimization of the policy and auxiliary losses. Importantly, the learning is done completely *on-policy*: the policy loss is computed from rollouts while the auxiliary losses use samples from a small replay buffer. Yet, even with these precautions, the authors are unable to leverage reconstruction by VAE (Kingma & Welling, 2013) and report its damaging affect on learning.

Similarly, Jaderberg et al. (2017) propose to use unsupervised auxiliary tasks, both observation and reward based, and show improvements in Atari, again in an *on-policy* regime<sup>2</sup>, which is much more stable for learning. Of all the auxiliary tasks considered by Jaderberg et al. (2017), reconstruction-based Pixel Control is the most effective. However, in maximizing changes in local patches, it imposes strong inductive biases that assume that dramatically changing pixel values and textures are correlated with good exploration and reward. Unfortunately, such highly task specific auxiliary is unlikely to scale to real world applications.

Generic pixel reconstruction is explored in Higgins et al. (2017b); Nair et al. (2018), where the authors use a beta variational autoencoder ( $\beta$ -VAE) (Kingma & Welling, 2013; Higgins et al., 2017a) and attempt to perform joint representation learning, but find it hard to train, thus reverting to the alternating training procedure (Lange & Riedmiller, 2010; Finn et al., 2015).

There has been more success in using model learning methods on images, such as Hafner et al. (2018); Lee et al. (2019). These methods use a world model approach (Ha & Schmidhuber, 2018), learning a representation space using a latent dynamics loss and pixel decoder loss to ground on the original observation space. These model-based reinforcement learning methods often show improved sample efficiency, but with the additional complexity of balancing various auxiliary losses, such as a dynamics loss, reward loss, and decoder loss in addition to the original policy and value optimizations. These proposed methods are correspondingly brittle to hyperparameter settings, and difficult to reproduce, as they balance multiple training objectives.

<sup>2</sup>Jaderberg et al. (2017) make use of a replay buffer that only stores the most recent 2K transitions, a small fraction of the 25M transitions experienced in training.

Task name	Number of Episodes	SAC:pixel	PlaNet	SLAC	SAC:state
finger_spin	1000	645 ± 37	659 ± 45	<b>900 ± 39</b>	<b>945 ± 19</b>
walker_walk	1000	33 ± 2	949 ± 9	864 ± 35	<b>974 ± 1</b>
ball_in_cup_catch	2000	593 ± 84	861 ± 80	932 ± 14	<b>981 ± 1</b>
cartpole.swingup	2000	758 ± 58	802 ± 19	-	<b>860 ± 8</b>
reacher_easy	2500	121 ± 28	<b>949 ± 25</b>	-	<b>953 ± 11</b>
cheetah_run	3000	366 ± 68	701 ± 6	<b>830 ± 32</b>	<b>836 ± 105</b>

Table 1. A comparison of current methods: SAC from pixels, PlaNet, SLAC, SAC from proprioceptive states (representing an upper bound). The large performance gap between SAC:pixel and SAC:state motivates us to address the representation learning bottleneck in model-free off-policy RL.

### 3. Background

#### 3.1. Markov Decision Process

A fully observable Markov decision process (MDP) can be described as  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $P(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$  is the transition probability distribution,  $R(\mathbf{s}_t, \mathbf{a}_t)$  is the reward function, and  $\gamma$  is the discount factor (Bellman, 1957). An agent starts in a initial state  $\mathbf{s}_1$  sampled from a fixed distribution  $p(\mathbf{s}_1)$ , then at each timestep  $t$  it takes an action  $\mathbf{a}_t \in \mathcal{A}$  from a state  $\mathbf{s}_t \in \mathcal{S}$  and moves to a next state  $\mathbf{s}_{t+1} \sim P(\cdot|\mathbf{s}_t, \mathbf{a}_t)$ . After each action the agent receives a reward  $r_t = R(\mathbf{s}_t, \mathbf{a}_t)$ . We consider episodic environments with the length fixed to  $T$ . The goal of standard RL is to learn a policy  $\pi(\mathbf{a}_t|\mathbf{s}_t)$  that can maximize the agent’s expected cumulative reward  $\sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi}[r_t]$ , where  $\rho_\pi$  is discounted state-action visitations of  $\pi$ , also known as occupancies. An important modification (Ziebart et al., 2008) augments this objective with an entropy term  $\mathcal{H}(\pi(\cdot|\mathbf{s}_t))$  to encourage exploration and robustness to noise. The resulting maximum entropy objective is then defined as

$$\pi^* = \arg \max_{\pi} \sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi}[r_t + \alpha \mathcal{H}(\pi(\cdot|\mathbf{s}_t))], \quad (1)$$

where  $\alpha$  is temperature that balances between optimizing for the reward and for the stochasticity of the policy.

#### 3.2. Soft Actor-Critic

Soft Actor-Critic (SAC) (Haarnoja et al., 2018) is an *off-policy* actor-critic method that uses the maximum entropy framework to derive soft policy iteration. At each iteration SAC performs soft policy evaluation and improvement steps. The policy evaluation step fits a parametric Q-function  $Q(\mathbf{s}_t, \mathbf{a}_t)$  using transitions sampled from the replay buffer  $\mathcal{D}$  by minimizing the soft Bellman residual

$$J(Q) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1}) \sim \mathcal{D}} \left[ (Q(\mathbf{s}_t, \mathbf{a}_t) - r_t - \gamma \bar{V}(\mathbf{s}_{t+1}))^2 \right]. \quad (2)$$

The target value function  $\bar{V}$  is approximated via a Monte-Carlo estimate of the following expectation

$$\bar{V}(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [\bar{Q}(\mathbf{s}_t, \mathbf{a}_t) - \alpha \log \pi(\mathbf{a}_t|\mathbf{s}_t)], \quad (3)$$

where  $\bar{Q}$  is the target Q-function parametrized by a weight vector obtained from an exponentially moving average of the Q-function weights to stabilize training. The policy improvement step then attempts to project a parametric policy  $\pi(\mathbf{a}_t|\mathbf{s}_t)$  by minimizing KL divergence between the policy and a Boltzmann distribution induced by the Q-function using the following objective

$$J(\pi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} [D_{\text{KL}}(\pi(\cdot|\mathbf{s}_t) || \mathcal{Q}(\mathbf{s}_t, \cdot))], \quad (4)$$

where  $\mathcal{Q}(\mathbf{s}_t, \cdot) \propto \exp\{\frac{1}{\alpha} Q(\mathbf{s}_t, \cdot)\}$ .

#### 3.3. Image-based Observations and Autoencoders

Directly learning from raw images posses an additional problem of partial observability, which is formalized by a partially observable MDP (POMDP). In this setting, instead of getting a low-dimensional state  $\mathbf{s}_t \in \mathcal{S}$  at time  $t$ , the agent receives a high-dimensional observation  $\mathbf{o}_t \in \mathcal{O}$ , which is a rendering of potentially incomplete view of the corresponding state  $\mathbf{s}_t$  of the environment (Kaelbling et al., 1998). This complicates applying RL as the agent now needs to also learn a compact latent representation to infer the state. Fitting a high-capacity encoder using only a scarce reward signal is sample inefficient and prone to suboptimal convergence. Following prior work (Lange & Riedmiller, 2010; Finn et al., 2015) we explore unsupervised pretraining via an image-based autoencoder (AE). In practice, the AE is represented as a convolutional encoder  $g_\phi$  that maps an image observation  $\mathbf{o}_t$  to a low-dimensional latent vector  $\mathbf{z}_t$ , and a deconvolutional decoder  $f_\theta$  that reconstructs  $\mathbf{z}_t$  back to the original image  $\mathbf{o}_t$ . Both the encoder and decoder are trained simultaneously by maximizing the expected log-likelihood

$$J(\text{AE}) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\log p_\theta(\mathbf{o}_t|\mathbf{z}_t)], \quad (5)$$

where  $\mathbf{z}_t = g_\phi(\mathbf{o}_t)$ . Or in the case of  $\beta$ -VAE (Kingma & Welling, 2013; Higgins et al., 2017a) we maximize the



Figure 1. Image-based continuous control tasks from the DeepMind Control Suite (Tassa et al., 2018) used in our experiments. Each task offers an unique set of challenges, including complex dynamics, sparse rewards, hard exploration, and other traits (see Appendix A).

objective below

$$J(\text{VAE}) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\mathbb{E}_{\mathbf{z}_t \sim q_\phi(\mathbf{z}_t | \mathbf{o}_t)} [\log p_\theta(\mathbf{o}_t | \mathbf{z}_t)] - \beta D_{\text{KL}}(q_\phi(\mathbf{z}_t | \mathbf{o}_t) || p(\mathbf{z}_t))], \quad (6)$$

where the variational distribution is parametrized as  $q_\phi(\mathbf{z}_t | \mathbf{o}_t) = \mathcal{N}(\mathbf{z}_t | \mu_\phi(\mathbf{o}_t), \sigma_\phi^2(\mathbf{o}_t))$ . The latent vector  $\mathbf{z}_t$  is then used by an RL algorithm, such as SAC, instead of the unavailable true state  $\mathbf{s}_t$ .

## 4. Representation Learning with Image Reconstruction

We start by noting a dramatic gap in an agent’s performance when it learns from image-based observations rather than low-dimensional proprioceptive states. Table 1 illustrates that in all cases **SAC:pixel** (an agent that learns from pixels) is significantly outperformed by **SAC:state** (an agent that learns from states). This result suggests that attaining a compact state representation is key in enabling efficient RL from images. Prior work has demonstrated that auxiliary supervision can improve representation learning, which is further confirmed in Table 1 by superior performance of model-based methods, such as PlaNet (Hafner et al., 2018) and SLAC (Lee et al., 2019), both of which make use of several auxiliary tasks to learn better representations.

While a wide range of auxiliary objectives could be added to aid effective representation learning, we focus our attention on the most general and widely applicable – an image reconstruction loss. Furthermore, our goal is to develop a simple and robust algorithm that has the potential to be scaled up to real world applications (e.g. robotics). Correspondingly, we avoid task dependent auxiliary losses, such as Pixel Control from Jaderberg et al. (2017), or world-models (Shelhamer et al., 2016; Hafner et al., 2018; Lee et al., 2019). As noted by Gelada et al. (2019) the latter can be brittle to train for reasons including: i) tension between reward and transition losses which requires careful tuning and ii) difficulty in modeling complex dynamics (which we explore further in Section 5.2).

Following Nair et al. (2018); Hafner et al. (2018); Lee et al. (2019), which use reconstruction loss to learn the representation space and dynamics model with a variational autoen-

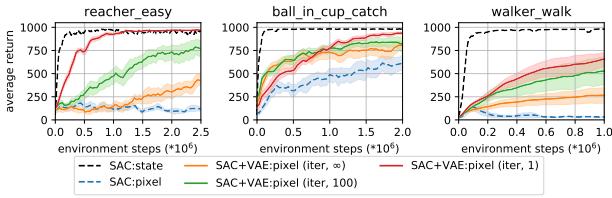
coder (Kingma & Welling, 2013; Higgins et al., 2017a), we also employ a  $\beta$ -VAE to learn representations, but in contrast to Hafner et al. (2018); Lee et al. (2019) we only consider reconstructing the current frame, instead of reconstructing a temporal sequence of frames. Based on evidence from Lange & Riedmiller (2010); Finn et al. (2015); Nair et al. (2018) we first try alternating between learning the policy and  $\beta$ -VAE, and in Section 4.2 observe a positive correlation between the alternation frequency and the agent’s performance. However, this approach does not fully close the performance gap, as the learned representation is not optimized for the task’s objective. To address this shortcoming, we then attempt to additionally update the  $\beta$ -VAE encoder with the actor-critic gradients. Unfortunately, our investigation in Section 4.3 shows this approach to be ineffective due to severe instability in training, especially with larger  $\beta$  values. Based on these results, in Section 4.4 we identify two reasons behind the instability, that originate from the stochastic nature of a  $\beta$ -VAE and the non-stationary gradient from the actor. We then propose two simple remedies and in Section 4.5 introduce our method for an effective model-free off-policy RL from images.

### 4.1. Experimental Setup

Before carrying out our empirical study, we detail the experimental setup. A more comprehensive overview can be found in Appendix B. We evaluate all agents on six challenging control tasks (Figure 1). For brevity, on occasion, results for three tasks are shown with the remainder presented in the appendix. An image observation is represented as a stack of three consecutive  $84 \times 84$  RGB renderings (Mnih et al., 2013) to infer temporal statistics, such as velocity and acceleration. For simplicity, we keep the hyper parameters fixed across all the tasks, except for action repeat (see Appendix B.3), which we set according to Hafner et al. (2018) for a fair comparison to the baselines. We evaluate an agent after every 10K training observations, by computing an average return over 10 episodes. For a reliable comparison we run 10 random seeds and report the mean and standard deviation of the evaluation reward.

## 4.2. Alternating Representation Learning with a $\beta$ -VAE

We first set out to confirm the benefits of an alternating approach to representation learning in off-policy RL. We conduct an experiment where we initially pretrain the convolutional encoder  $g_\phi$  and deconvolutional decoder  $f_\theta$  of a  $\beta$ -VAE according to the loss  $J(\text{VAE})$  (Equation (6)) on observations collected by a random policy. The actor and critic networks of SAC are then trained for  $N$  steps using latent states  $\mathbf{z}_t \sim g_\phi(\mathbf{o}_t)$  as inputs instead of image-based observations  $\mathbf{o}_t$ . We keep the encoder  $g_\phi$  fixed during this period. The updated policy is then used to interact with the environment to gather new transitions that are consequently stored in the replay buffer. We continue iterating between the autoencoder and actor-critic updates until convergence. Note that the gradients are never shared between the  $\beta$ -VAE for learning the representation space, and the actor-critic. In Figure 2 we vary the frequency  $N$  at which the representation space is updated, from  $N = \infty$  where the representation is never updated after the initial pretraining period, to  $N = 1$  where the representation is updated after every policy update. We observe a positive correlation between this frequency and the agent's performance. Although the alternating scheme helps to improve the sample efficiency of the agent, it still falls short of reaching the upper bound performance of SAC:state. This is not surprising, as the learned representation space is never optimized for the task's objective.



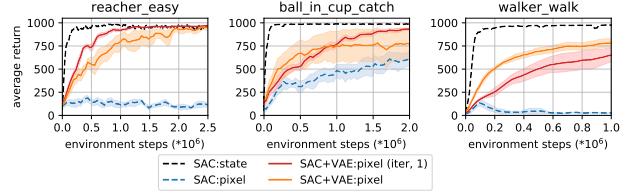
**Figure 2.** Separate  $\beta$ -VAE and policy training with no shared gradients SAC+VAE:pixel ( $\text{iter. } N$ ), with SAC:state shown as an upper bound.  $N$  refers to frequency in environment steps at which the  $\beta$ -VAE updates after initial pretraining. More frequent updates are beneficial for learning better representations, but cannot fully address the gap in performance. Full results in Appendix C.

## 4.3. Joint Representation Learning with a $\beta$ -VAE

To further improve performance of the agent we seek to learn a latent representation that is well aligned with the underlying RL objective. Shelhamer et al. (2016) has demonstrated that joint policy and auxiliary objective optimization improves on the pretraining approach, as described in Section 4.2, but this has been only shown in the *on-policy* regime.

Thus we now attempt to verify the feasibility of joint representation learning with a  $\beta$ -VAE in the *off-policy* setting.

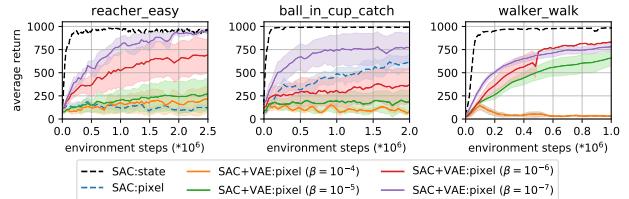
Specifically, we want to update the encoder network  $g_\phi$  with the gradients coming through the latent state  $\mathbf{z}_t$  from the actor  $J(\pi)$  (Equation (4)), critic  $J(Q)$  (Equation (2)), and  $\beta$ -VAE  $J(\text{VAE})$  (Equation (6)) losses. We thus take the best performing variant from the previous experiment (e.g. SAC+VAE:pixel (iter, 1)) and let the actor-critic's gradients update the encoder  $g_\phi$ . We tune for the best  $\beta$  and name this agent **SAC+VAE:pixel**. Results in Figure 3 show that the joint representation learning with  $\beta$ -VAE in unstable in the *off-policy* setting and performs worse than the baseline that does not utilize task dependent information (e.g. SAC+VAE:pixel (iter, 1)).



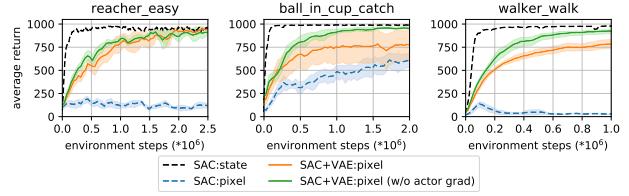
**Figure 3.** An unsuccessful attempt to propagate gradients from the actor-critic down to the  $\beta$ -VAE encoder. SAC+VAE:pixel exhibits instability in training which leads to subpar performance comparing to the baseline SAC+VAE:pixel (iter, 1), which does not use the actor-critic gradients. Full results in Appendix D.

## 4.4. Stabilizing Joint Representation Learning

Following an unsuccessful attempt at joint representation learning with a  $\beta$ -VAE in off-policy RL, we investigate the root cause of the instability.



(a) Smaller values of  $\beta$  reduce stochasticity of a  $\beta$ -VAE and lead to a better performance.



(b) Preventing the actor's gradients to update the convolutional encoder helps to improve performance even further.

**Figure 4.** We identify two reasons for the subpar performance of joint representation learning. (a) The stochastic nature of a  $\beta$ -VAE, and (b) the non-stationary actor's gradients. Full results in Appendix E.

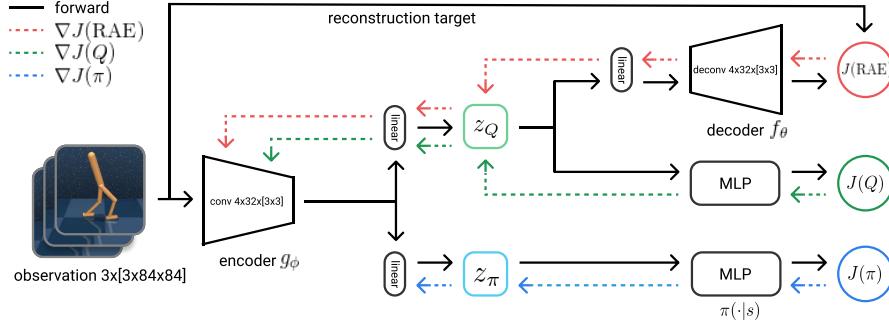


Figure 5. Our algorithm (SAC+AE) augments SAC with a regularized autoencoder to achieve stable training from images in the off-policy regime. The stability comes from switching to a deterministic encoder that is carefully updated with gradients from the reconstruction  $J(\text{RAE})$  (Equation (7)) and soft Q-learning  $J(Q)$  (Equation (2)) objectives.

We first observe that the stochastic nature of a  $\beta$ -VAE damages performance of the agent. The results from Figure 4a illustrate that smaller values of  $\beta$  improve the training stability as well as the task performance. This motivates us to instead consider a completely deterministic autoencoder.

Furthermore, we observe that updating the convolutional encoder with the actor’s gradients hurts the agent’s performance. In Figure 4b we observe that blocking the actor’s gradients from propagating to the encoder improves results considerably. This is because updating the encoder with the  $J(\pi)$  loss (Equation (4)) also changes the  $Q$ -function network inside the objective, due to the convolutional encoder being shared between the policy  $\pi$  and  $Q$ -function. A similar phenomenon has been observed by Mnih et al. (2013), where the authors employ a static target  $Q$ -function to stabilize TD learning. It might appear that updating the encoder with only the critic’s gradients would be insufficient to learn a task dependent representation space. However, the policy  $\pi$  in SAC is a parametric projection of a Boltzmann distribution induced by the  $Q$ -function, see Equation (4). Thus, the  $Q$ -function contains all relevant information about the task and allows the encoder to learn task dependent representations from the critic’s gradient alone.

#### 4.5. Our Approach SAC+AE: Joint Off-Policy Representation Learning

We now introduce our approach **SAC+AE** – a stable off-policy RL algorithm from images, derived from the above findings. We first replace the  $\beta$ -VAE with a deterministic autoencoder. To preserve the regularization effects of a  $\beta$ -VAE we adopt the RAE approach of Ghosh et al. (2019), which imposes a  $L_2$  penalty on the learned representation  $\mathbf{z}_t$  and weight-decay on the decoder parameters

$$J(\text{RAE}) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\log p_\theta(\mathbf{o}_t | \mathbf{z}_t) + \lambda_z \|\mathbf{z}_t\|^2 + \lambda_\theta \|\theta\|^2], \quad (7)$$

where  $\mathbf{z}_t = g_\phi(\mathbf{o}_t)$ , and  $\lambda_z$ ,  $\lambda_\theta$  are hyper parameters.

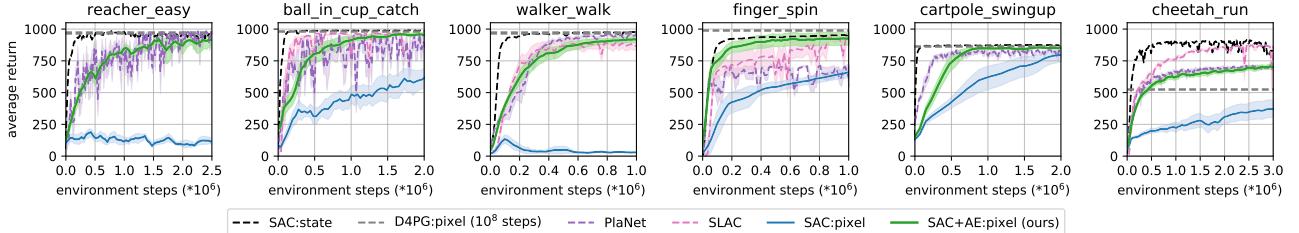
We also prevent the actor’s gradients from updating the convolutional encoder, as suggested in Section 4.4. Unfortunately, this slows down signal propagation to the encoder, and thus we find it important to update the convolutional weights of the target  $Q$ -function faster than the rest of the network’s parameters. We thus employ different rates  $\tau_Q$  and  $\tau_{\text{enc}}$  (with  $\tau_{\text{enc}} > \tau_Q$ ) to compute Polyak averaging over the corresponding parameters of the target  $Q$ -function. Our approach is summarized in Figure 5.

## 5. Evaluation of SAC+AE

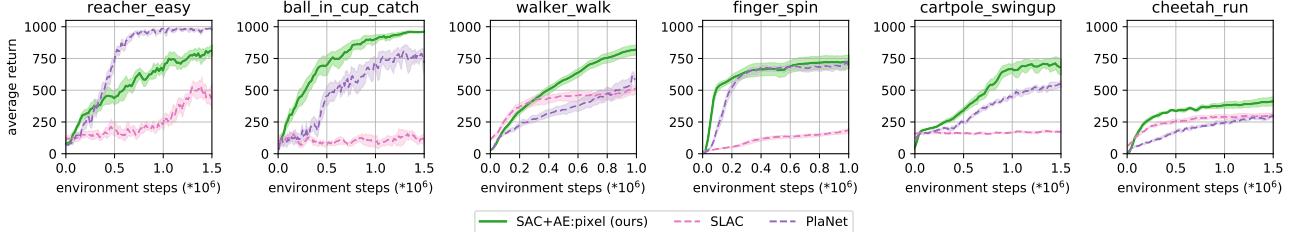
In this section we evaluate our approach, **SAC+AE**, on various benchmark tasks and compare against state-of-the-art methods, both model-free and model-based. We then highlight the benefits of our model-free approach over those model-based methods in modified environments with distractors, as an approximation of real world noise. Finally, we test generalization to unseen tasks and dissect the representation power of the encoder.

### 5.1. Learning Control from Pixels

We evaluate our method on six challenging image-based continuous control tasks (see Figure 1) from DMC (Tassa et al., 2018) and compare against several state-of-the-art model-free and model-based RL algorithms for learning from pixels: **D4PG** (Barth-Maron et al., 2018), an off-policy actor-critic algorithm; **PlaNet** (Hafner et al., 2018), a model-based method that learns a dynamics model with deterministic and stochastic latent variables and employs cross-entropy planning for control; and **SLAC** (Lee et al., 2019), which combines a purely stochastic latent model together with an model-free soft actor-critic. In addition, we compare against SAC:state that learns from low-dimensional proprioceptive state, as an upper bound on performance. Results in Figure 6a illustrate that **SAC+AE:pixel** matches the state-of-the-art model-based methods such as PlaNet and SLAC, despite being



(a) Our method demonstrates significantly improved performance over the baseline SAC:pixel. Moreover, it matches the state-of-the-art performance of world-model based algorithms, such as PlaNet and SLAC, as well as a model-free algorithm D4PG, that learns directly from raw images. Our algorithm is extremely stable, robust, and straightforward to implement.



(b) Methods that rely on forward modeling, such as PlaNet and SLAC, suffer severely from the background noise, while our approach is resistant to the distractors. Examples of background distractors are shown in Figure 7.

Figure 6. Two main results of our work. In (a) we demonstrate that our simple method matches the state-of-the-art performance on DMC tasks. In (b) we outperform the baselines on more complicated tasks where the observations are altered with noise.

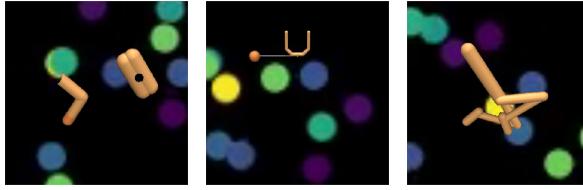


Figure 7. Backgrounds altered with randomly moving distractors.

extremely simple and straightforward to implement.

## 5.2. Performance on Noisy Observations

Performing accurate forward-modeling predictions based off of noisy observations is challenging and requires learning a high fidelity model that encapsulates strong inductive biases (Watters et al., 2017). The current state-of-the-art world-model based approaches (Hafner et al., 2018; Lee et al., 2019) solely rely on a general purpose recurrent state-space model parametrized with a  $\beta$ -VAE, and thus are highly vulnerable to the observational noise. In contrast, the representations learned with just reconstruction loss are better suited to handle the background noise.

To confirm this, we evaluate several agents on tasks where we add simple distractors in the background, consisting of colored balls bouncing off each other and the frame (Figure 7). We use image processing to filter away the static background and replace it with this dynamic noise, as proposed in Zhang et al. (2018b). We aim to emulate a common setup in a robotics lab, where various unrelated objects can affect robot’s observations. In Figure 6b we see

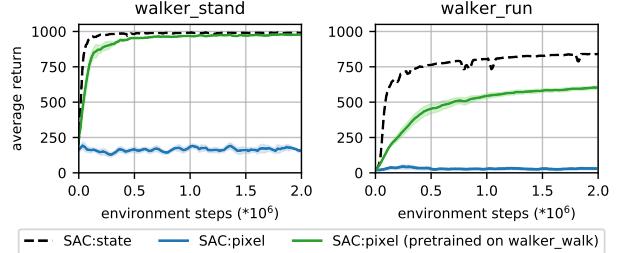
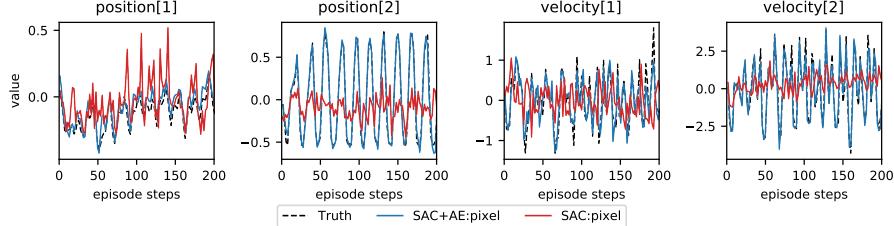


Figure 8. An encoder pretrained with our method (SAC+AE:pixel) on walker\_walk is able to generalize to unseen walker\_stand and walker\_run tasks. All three tasks share similar image observations, but have quite different reward structure. SAC with a pretrained on walker\_walk encoder significantly outperforms the baseline.

that methods that rely on forward modeling perform drastically worse than our approach, showing that our method is more robust to background noise.

## 5.3. Generalization to Unseen Tasks

Next, we show that the latent representation space learned by our method is able to generalize to different tasks without additional fine-tuning. We take three tasks walker\_stand, walker\_walk, and walker\_run from DMC, which share the same observational appearance, but all have different reward functions. We train SAC+AE:pixel on the walker\_walk task until convergence and fix the encoder. Consequently, we train two



**Figure 9.** Linear projections of latent representation spaces learned by our method (SAC+AE:pixel) and the baseline (SAC:pixel) onto proprioceptive states. We compare ground truth value of each proprioceptive coordinate against their reconstructions for `cheetah_run`, and conclude that our method successfully encodes proprioceptive state information. For visual clarity we only plot 2 position (out of 8) and 2 velocity (out of 9) coordinates. Full results in [Appendix G](#).

SAC:pixel agents on `walker_stand` and `walker_run`. The encoder of the first agent is initialized with weights from the pretrained on `walker_walk` encoder, while the encoder of the second agent is not. Neither of the agents uses the reconstruction signal, and only backpropagate the critic’s gradients. Results in [Figure 8](#) illustrate that our method learns latent representations that can readily generalize to unseen tasks and achieve much better performance than SAC:pixel trained from scratch.

#### 5.4. Representation Power of the Encoder

Finally, we want to determine if our method is able to extract sufficient information from raw images to recover the corresponding proprioceptive states. We thus train SAC+AE:pixel and SAC:pixel until convergence on `cheetah_run` and then fix their encoders. We then learn two linear projections to map the encoders’ latent embedding of image observations into the corresponding proprioceptive states. Finally, we compare ground truth proprioceptive states against their reconstructions. We emphasize that the image encoder attributes for over 90% of the agent’s parameters, thus we believe that the encoder’s latent output  $\mathbf{z}_t$  captures a significant amount of information about the corresponding internal state in both cases, even though SAC:pixel does not require this explicitly. Results in [Figure 9](#) confirm that the internals of the task are easily extracted from the encoder grounded on pixel observations, whereas they are much more difficult to construct from the representation learned by SAC:pixel.

## 6. Discussion

For RL agents to be effective in the real world, where vision is one of the richest sensing modalities, we need sample efficient, robust algorithms that work from pixel observations. We pinpoint two strategies to obtain sample efficiency – i) use off-policy methods and ii) use self-supervised auxiliary losses. For methods to be robust, we want auxiliary losses that do not rely on task-specific inductive biases, so we focus on a simple reconstruction loss. In this work, we provide a thorough study into combining

reconstruction loss with off-policy methods for improved sample efficiency in rich observation settings. Our analysis yields two key findings. The first is that deterministic AE models outperform  $\beta$ -VAEs ([Higgins et al., 2017a](#)), due to additional instabilities such as bootstrapping, off-policy data, and joint training with auxiliary losses. The second is that propagating the actor’s gradients through the convolutional encoder hurts performance.

Based on these results, we also recommend an effective off-policy, model-free RL algorithm for pixel observations with only reconstruction loss as an auxiliary task. It is competitive with state-of-the-art model-based methods on traditional benchmarks, but much simpler, robust, and does not require learning a dynamics model ([Figure 6a](#)). We show through ablations the superiority of joint learning over previous methods that use an alternating training procedure with separated gradients, the necessity of a pixel reconstruction loss over reconstruction to lower-dimensional “correct” representations, and demonstrations of the representation power and generalization ability of our learned representation. We additionally construct settings with distractors approximating real world noise which show how learning a world-model as an auxiliary loss can be harmful ([Figure 6b](#)), and in which our method, SAC+AE, exhibits state-of-the-art performance.

In the Appendix we provide results across all experiments on the full suite of 6 tasks chosen from DMC ([Appendix A](#)), and the full set of hyperparameters used in [Appendix B](#). There are also additional experiments autoencoder capacity ([Appendix F](#)), a look at optimality of the learned latent representation ([Appendix I](#)) and importance of action repeat ([Appendix J](#)). Finally, we opensource our codebase for the community to spur future research in image-based RL.

## References

- Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. *arXiv e-prints*, 2016.
- Barth-Maron, G., Hoffman, M. W., Budden, D., Dabney, W., Horgan, D., TB, D., Muldal, A., Heess, N., and Lillicrap, T. Distributional policy gradients. In *International Conference on Learning Representations*, 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*, 2013.
- Bellman, R. A markovian decision process. *Indiana Univ. Math. J.*, 1957.
- Dwibedi, D., Tompson, J., Lynch, C., and Sermanet, P. Learning actionable representations from visual observations. *CoRR*, 2018.
- Finn, C., Tan, X. Y., Duan, Y., Darrell, T., Levine, S., and Abbeel, P. Learning visual feature spaces for robotic manipulation with deep spatial autoencoders. *CoRR*, 2015.
- Gelada, C., Kumar, S., Buckman, J., Nachum, O., and Bellemare, M. G. Deepmdp: Learning continuous latent space models for representation learning. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Ghosh, P., Sajjadi, M. S. M., Vergari, A., Black, M. J., and Schölkopf, B. From variational to deterministic autoencoders. *arXiv e-prints*, 2019.
- Ha, D. and Schmidhuber, J. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., 2018.
- Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent dynamics for planning from pixels. *arXiv preprint arXiv:1811.04551*, 2018.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. betavae: Learning basic visual concepts with a constrained variational framework. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017a.
- Higgins, I., Pal, A., Rusu, A. A., Matthey, L., Burgess, C. P., Pritzel, A., Botvinick, M., Blundell, C., and Lerchner, A. Darla: Improving zero-shot transfer in reinforcement learning. *CoRR*, 2017b.
- Jaderberg, M., Mnih, V., Czarnecki, W., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. *International Conference on Learning Representations*, 2017.
- Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. *arXiv e-prints*, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Lange, S. and Riedmiller, M. A. Deep auto-encoder neural networks in reinforcement learning. In *IJCNN*, 2010.
- Lee, A. X., Nagabandi, A., Abbeel, P., and Levine, S. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. *arXiv e-prints*, 2019.
- Mattner, J., Lange, S., and Riedmiller, M. Learn to swing up and balance a real pole based on raw visual input data. In *Neural Information Processing*, 2012.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv e-prints*, 2013.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T. P., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. *CoRR*, 2016.
- Munk, J., Kober, J., and Babuska, R. Learning state representation for deep actor-critic control. In *Proceedings 2016 IEEE 55th Conference on Decision and Control (CDC)*, 2016.
- Nair, A. V., Pong, V., Dalal, M., Bahl, S., Lin, S., and Levine, S. Visual reinforcement learning with imagined goals. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 9191–9200. Curran Associates, Inc., 2018.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv e-prints*, 2013.

- Shelhamer, E., Mahmoudieh, P., Argus, M., and Darrell, T.  
Loss is its own reward: Self-supervision for reinforcement learning. *CoRR*, 2016.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- van Hasselt, H., Guez, A., and Silver, D. Deep reinforcement learning with double q-learning. *arXiv e-prints*, 2015.
- Watters, N., Tacchetti, A., Weber, T., Pascanu, R., Battaglia, P. W., and Zoran, D. Visual interaction networks. *CoRR*, 2017.
- Xiao, L., Bahri, Y., Sohl-Dickstein, J., Schoenholz, S. S., and Pennington, J. Dynamical isometry and a mean field theory of cnns: How to train 10,000-layer vanilla convolutional neural networks. *arXiv e-prints*, 2018.
- Yarats, D. and Kostrikov, I. Soft actor-critic (sac) implementation in pytorch. [https://github.com/denisyarats/pytorch\\_sac](https://github.com/denisyarats/pytorch_sac), 2020.
- Zhang, A., Satija, H., and Pineau, J. Decoupling dynamics and reward for transfer learning. *CoRR*, 2018a.
- Zhang, A., Wu, Y., and Pineau, J. Natural environment benchmarks for reinforcement learning. *CoRR*, 2018b.
- Ziebart, B. D., Maas, A., Bagnell, J. A., and Dey, A. K. Maximum entropy inverse reinforcement learning. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3*, 2008.

## Appendix

### A. The DeepMind Control Suite

We evaluate the algorithms in the paper on the DeepMind control suite (DMC) (Tassa et al., 2018) – a collection of continuous control tasks that offers an excellent testbed for reinforcement learning agents. The software emphasizes the importance of having a standardised set of benchmarks with a unified reward structure in order to measure made progress reliably.

Specifically, we consider six domains (see Figure 10) that result in twelve different control tasks. Each task (Table 2) poses a particular set of challenges to a learning algorithm. The `ball_in_cup_catch` task only provides the agent with a sparse reward when the ball is caught; the `cheetah_run` task offers high dimensional internal state and action spaces; the `reacher_hard` task requires the agent to explore the environment. We refer the reader to the original paper to find more information about the benchmarks.

Task name	dim( $\mathcal{O}$ )		dim( $\mathcal{A}$ )	Reward type
	Proprioceptive	Image-based		
<code>ball_in_cup_catch</code>	8	$3 \times 84 \times 84$	2	sparse
<code>cartpole_{balance, swingup}</code>	5	$3 \times 84 \times 84$	1	dense
<code>cheetah_run</code>	17	$3 \times 84 \times 84$	6	dense
<code>finger_{spin, turn_easy, turn_hard}</code>	12	$3 \times 84 \times 84$	2	dense/sparse
<code>reacher_{easy, hard}</code>	7	$3 \times 84 \times 84$	2	sparse
<code>walker_{stand, walk, run}</code>	24	$3 \times 84 \times 84$	6	dense

Table 2. Specifications of observation space  $\mathcal{O}$  (proprioceptive and image-based), action space  $\mathcal{A}$ , and the reward type for each task.

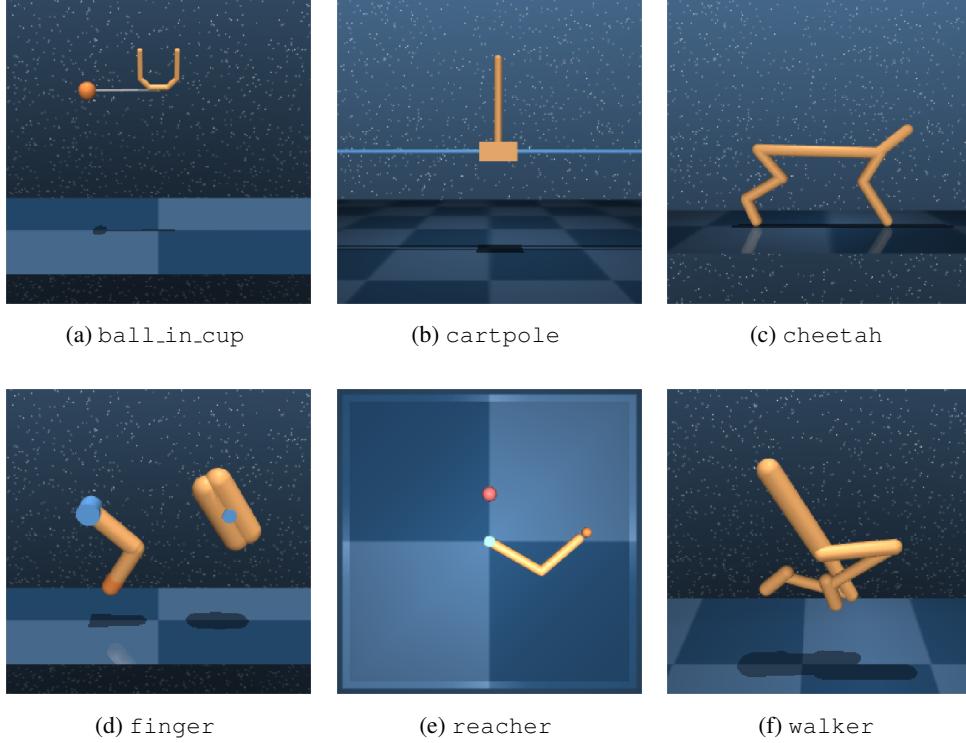


Figure 10. Our testbed consists of six domains spanning the total of twelve challenging continuous control tasks: `finger_{spin, turn_easy, turn_hard}`, `cartpole_{balance, swingup}`, `cheetah_run`, `walker_{stand, walk, run}`, `reacher_{easy, hard}`, and `ball_in_cup_catch`.

## B. Hyper Parameters and Setup

Our PyTorch SAC ([Haarnoja et al., 2018](#)) implementation is based off of ([Yarats & Kostrikov, 2020](#)).

### B.1. Actor and Critic Networks

We employ double Q-learning ([van Hasselt et al., 2015](#)) for the critic, where each Q-function is parametrized as a 3-layer MLP with ReLU activations after each layer except of the last. The actor is also a 3-layer MLP with ReLUs that outputs mean and covariance for the diagonal Gaussian that represents the policy. The hidden dimension is set to 1024 for both the critic and actor.

### B.2. Encoder and Decoder Networks

We employ an almost identical encoder architecture as in [Tassa et al. \(2018\)](#), with two minor differences. Firstly, we add two more convolutional layers to the convnet trunk. Secondly, we use ReLU activations after each conv layer, instead of ELU. We employ kernels of size  $3 \times 3$  with 32 channels for all the conv layers and set stride to 1 everywhere, except of the first conv layer, which has stride 2. We then take the output of the convnet and feed it into a single fully-connected layer normalized by LayerNorm ([Ba et al., 2016](#)). Finally, we add  $\tanh$  nonlinearity to the 50 dimensional output of the fully-connected layer.

The actor and critic networks both have separate encoders, although we share the weights of the conv layers between them. Furthermore, only the critic optimizer is allowed to update these weights (e.g. we truncate the gradients from the actor before they propagate to the shared conv layers).

The decoder consists of one fully-connected layer that is then followed by four deconv layers. We use ReLU activations after each layer, except the final deconv layer that produces pixels representation. Each deconv layer has kernels of size  $3 \times 3$  with 32 channels and stride 1, except of the last layer, where stride is 2.

We then combine the critic’s encoder together with the decoder specified above into an autoencoder. Note, because we share conv weights between the critic’s and actor’s encoders, the conv layers of the actor’s encoder will be also affected by reconstruction signal from the autoencoder.

### B.3. Training and Evaluation Setup

We first collect 1000 seed observations using a random policy. We then collect training observations by sampling actions from the current policy. We perform one training update every time we receive a new observation. In cases where we use action repeat, the number of training observations is only a fraction of the environment steps (e.g. a 1000 steps episode at action repeat 4 will only results into 250 training observations). The action repeat used for each environment is specified in [Table 3](#), following those used by PlaNet and SLAC.

We evaluate our agent after every 10000 environment steps by computing an average episode return over 10 evaluation episodes. Instead of sampling from the Gaussian policy we take its mean during evaluation.

We preserve this setup throughout all the experiments in the paper.

Task name	Action repeat
cartpole_swingup	8
reacher_easy	4
cheetah_run	4
finger_spin	2
ball_in_cupCatch	4
walker_walk	2

*Table 3.* Action repeat parameter used per task, following PlaNet and SLAC.

#### B.4. Weights Initialization

We initialize the weight matrix of fully-connected layers with the orthogonal initialization (Saxe et al., 2013) and set the bias to be zero. For convolutional and deconvolutional layers we use delta-orthogonal initialization (Xiao et al., 2018).

#### B.5. Regularization

We regularize the autoencoder network using the scheme proposed in Ghosh et al. (2019). In particular, we extend the standard reconstruction loss for a deterministic autoencoder with a  $L_2$  penalty on the learned representation  $\mathbf{z}$  and add weight decay on the decoder parameters  $\theta$ :

$$J(\text{RAE}) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\log p_\theta(\mathbf{o}_t | \mathbf{z}_t) + \lambda_z \|\mathbf{z}_t\|^2 + \lambda_\theta \|\theta\|^2] \quad \text{where} \quad \mathbf{z}_t = g_\phi(\mathbf{o}_t). \quad (8)$$

We set  $\lambda_z = 10^{-6}$  and  $\lambda_\theta = 10^{-7}$ .

#### B.6. Pixels Preprocessing

We construct an observational input as an 3-stack of consecutive frames (Mnih et al., 2013), where each frame is a RGB rendering of size  $84 \times 84$  from the 0th camera. We then divide each pixel by 255 to scale it down to  $[0, 1]$  range. For reconstruction targets we instead preprocess images by reducing bit depth to 5 bits as in Kingma & Dhariwal (2018).

#### B.7. Other Hyper Parameters

We also provide a comprehensive overview of all the remaining hyper parameters in Table 4.

Parameter name	Value
Replay buffer capacity	1000000
Batch size	128
Discount $\gamma$	0.99
Optimizer	Adam
Critic learning rate	$10^{-3}$
Critic target update frequency	2
Critic Q-function soft-update rate $\tau_Q$	0.01
Critic encoder soft-update rate $\tau_{\text{enc}}$	0.05
Actor learning rate	$10^{-3}$
Actor update frequency	2
Actor log stddev bounds	$[-10, 2]$
Autoencoder learning rate	$10^{-3}$
Temperature learning rate	$10^{-4}$
Temperature Adam's $\beta_1$	0.5
Init temperature	0.1

Table 4. A complete overview of used hyper parameters.

### C. Alternating Representation Learning with a $\beta$ -VAE

Iterative pretraining suggested in Lange & Riedmiller (2010); Finn et al. (2015) allows for faster representation learning, which consequently boosts the final performance, yet it is not sufficient enough to fully close the gap and additional modifications, such as joint training, are needed. Figure 11 provides additional results for the experiment described in Section 4.2.

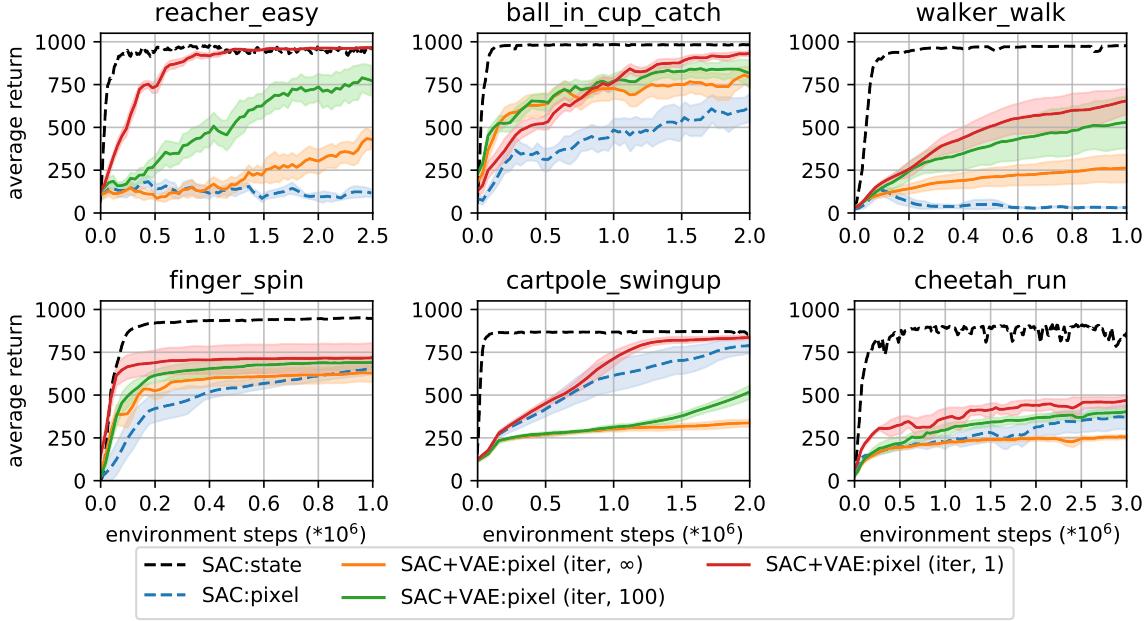
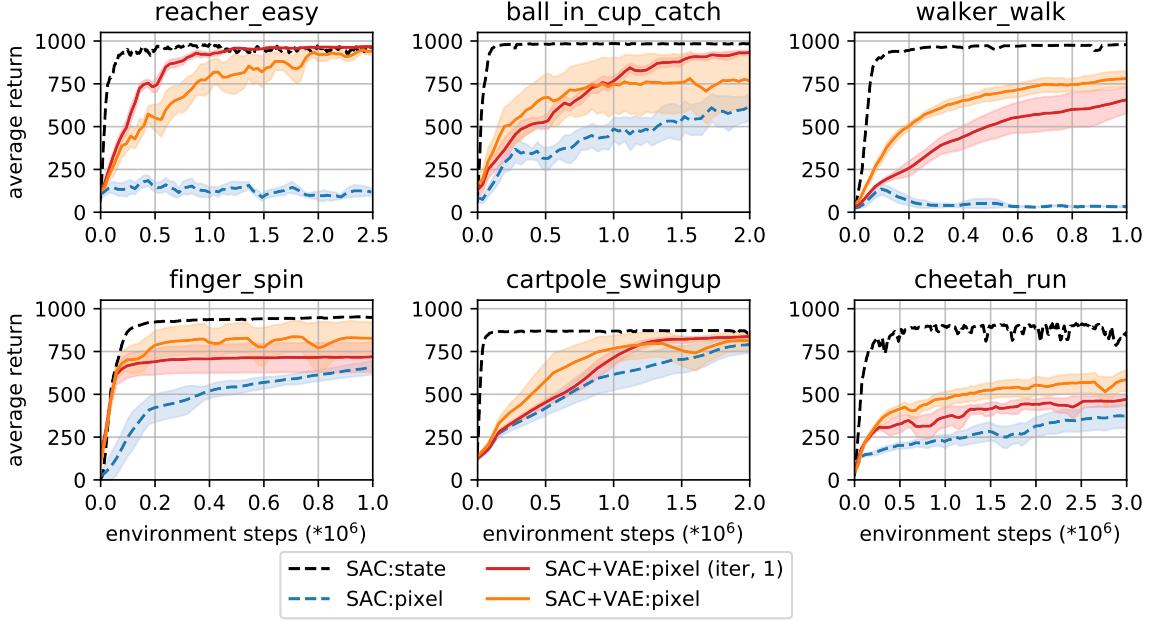


Figure 11. Separate  $\beta$ -VAE and policy training with no shared gradients SAC+VAE:pixel ( $\text{iter}, N$ ), with SAC:state shown as an upper bound.  $N$  refers to frequency in environment steps at which the  $\beta$ -VAE updates after initial pretraining. More frequent updates are beneficial for learning better representations, but cannot fully address the gap in performance.

## D. Joint Representation Learning with a $\beta$ -VAE

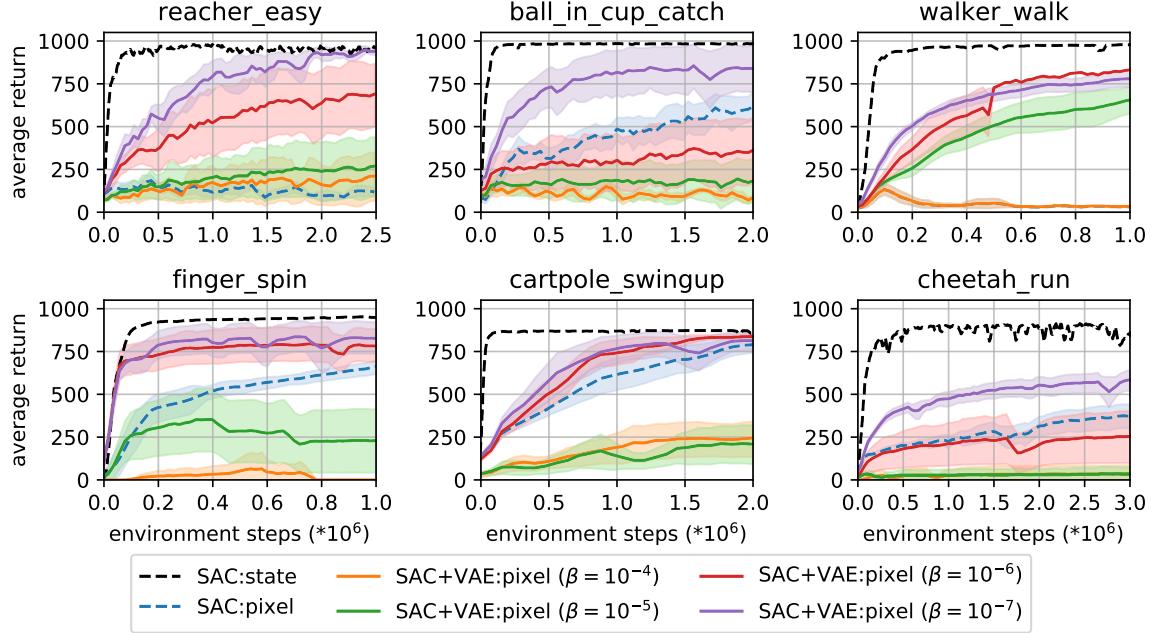
Additional results to the experiments from Section 4.3 are in Figure 4a and Figure 12.



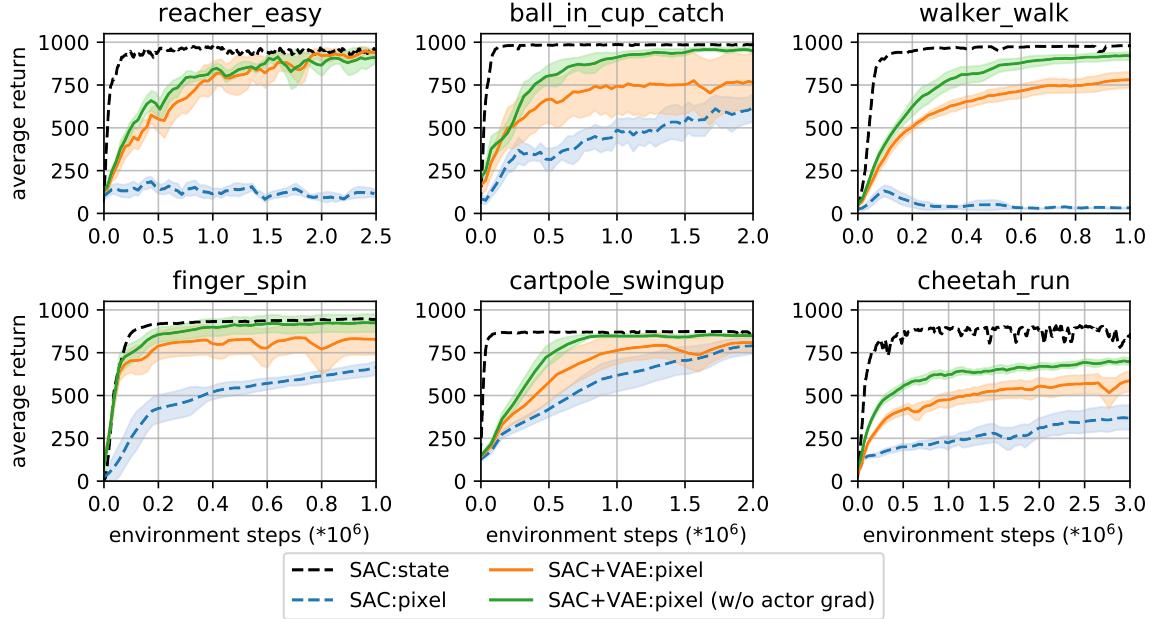
*Figure 12.* An unsuccessful attempt to propagate gradients from the actor-critic down to the encoder of the  $\beta$ -VAE to enable joint off-policy training. The learning process of SAC+VAE:pixel exhibits instability together with the subpar performance comparing to the baseline SAC+VAE:pixel (iter, 1), which does not share gradients with the actor-critic.

## E. Stabilizing Joint Representation Learning

Additional results to the experiments from Section 4.4 are in Figure 13.



(a) Smaller values of  $\beta$  reduce stochasticity of a  $\beta$ -VAE and lead to a better performance.

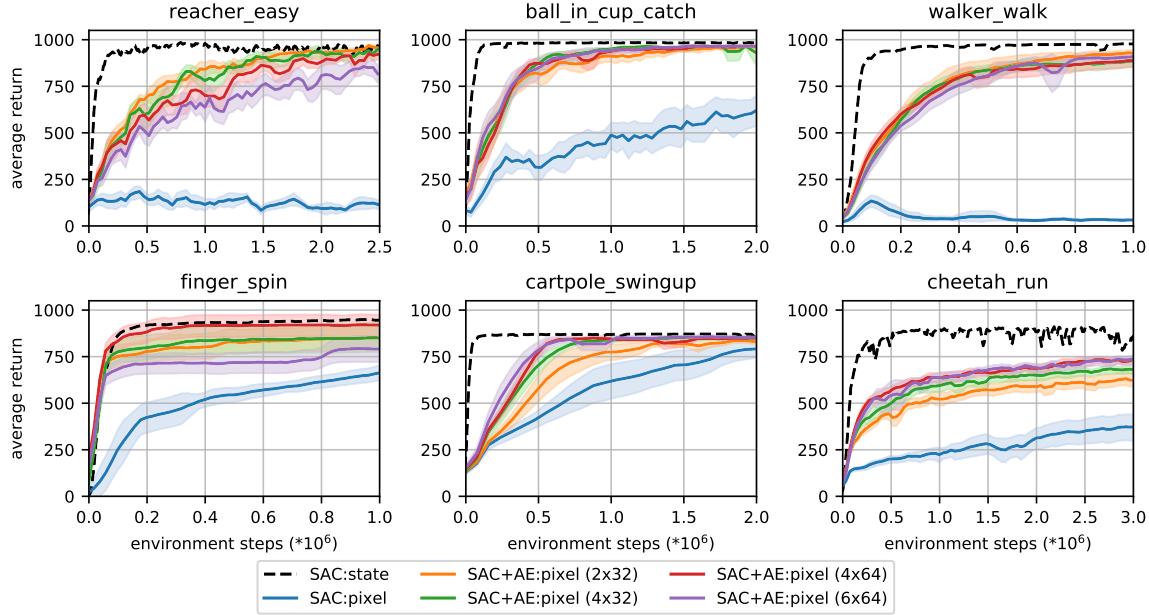


(b) Preventing the actor's gradients to update the convolutional encoder helps to improve performance even further.

Figure 13. We identify two reasons for the subpar performance of joint representation learning. (a) The stochastic nature of a  $\beta$ -VAE, and (b) the non-stationary actor's gradients.

## F. Capacity of the Autoencoder

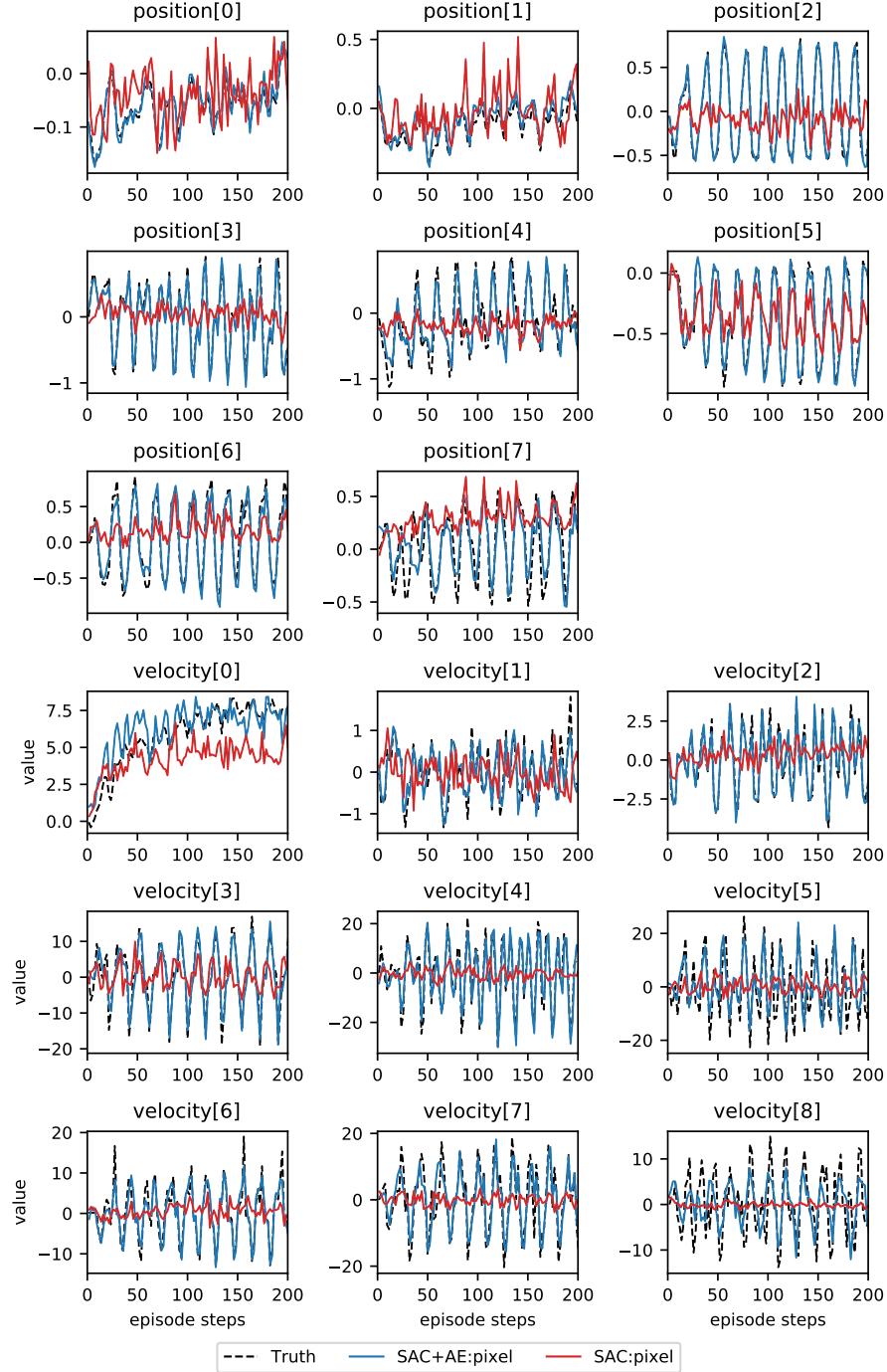
We also investigate various autoencoder capacities for the different tasks. Specifically, we measure the impact of changing the capacity of the convolutional trunk of the encoder and corresponding deconvolutional trunk of the decoder. Here, we maintain the shared weights across convolutional layers between the actor and critic, but modify the number of convolutional layers and number of filters per layer in Figure 14 across several environments. We find that SAC+AE is robust to various autoencoder capacities, and all architectures tried were capable of extracting the relevant features from pixel space necessary to learn a good policy. We use the same training and evaluation setup as detailed in Appendix B.3.



*Figure 14.* Different autoencoder architectures, where we vary the number of conv layers and the number of output channels in each layer in both the encoder and decoder. For example,  $4 \times 32$  specifies an architecture with 4 conv layers, each outputting 32 channels. We observe that the difference in capacity has only limited effect on final performance.

## G. Representation Power of the Encoder

Addition results to the experiment in Section 5.4 that demonstrates encoder’s power to reconstruct proprioceptive state from image-observations are shown in Figure 15.



*Figure 15.* Linear projections of latent representation spaces learned by our method (SAC+AE:pixel) and the baseline (SAC:pixel) onto proprioceptive states. We compare ground truth value of each proprioceptive coordinate against their reconstructions for `cheetah_run`, and conclude that our method successfully encodes proprioceptive state information. The proprioceptive state of `cheetah_run` has 8 position and 9 velocity coordinates.

## H. Decoding to Proprioceptive State

Learning from low-dimensional proprioceptive observations achieves better final performance with greater sample efficiency (see Figure 6a for comparison to pixels baselines), therefore our intuition is to directly use these compact observations as the reconstruction targets to generate an auxiliary signal. Although, this is an unrealistic setup, given that we do not have access to proprioceptive states in practice, we use it as a tool to understand if such supervision is beneficial for representation learning and therefore can achieve good performance. We augment the observational encoder  $g_\phi$ , that maps an image  $\mathbf{o}_t$  into a latent vector  $\mathbf{z}_t$ , with a state decoder  $f_\theta$ , that restores the corresponding state  $\mathbf{s}_t$  from the latent vector  $\mathbf{z}_t$ . This leads to an auxiliary objective  $\mathbb{E}_{\mathbf{o}_t, \mathbf{s}_t \sim \mathcal{D}} [\frac{1}{2} \|f_\theta(\mathbf{z}_t) - \mathbf{s}_t\|_2^2]$ , where  $\mathbf{z}_t = g_\phi(\mathbf{o}_t)$ . We parametrize the state decoder  $f_\theta$  as a 3-layer MLP with 1024 hidden size and ReLU activations, and train it jointly with the actor-critic network. Such auxiliary supervision helps less than expected, and surprisingly hurts performance in `ball_in_cup_catch`, as seen in Figure 16. Our intuition is that such low-dimensional supervision is not able to provide the rich reconstruction error needed to fit the high-capacity convolutional encoder  $g_\phi$ . We thus seek for a denser auxiliary signal and try learning latent representation spaces with pixel reconstructions.

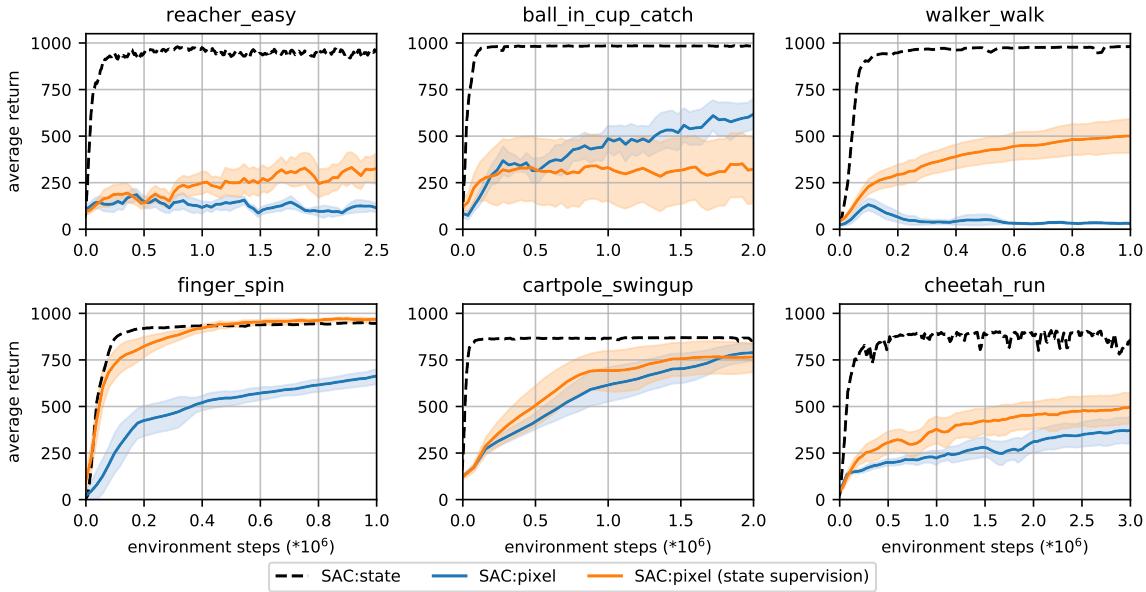
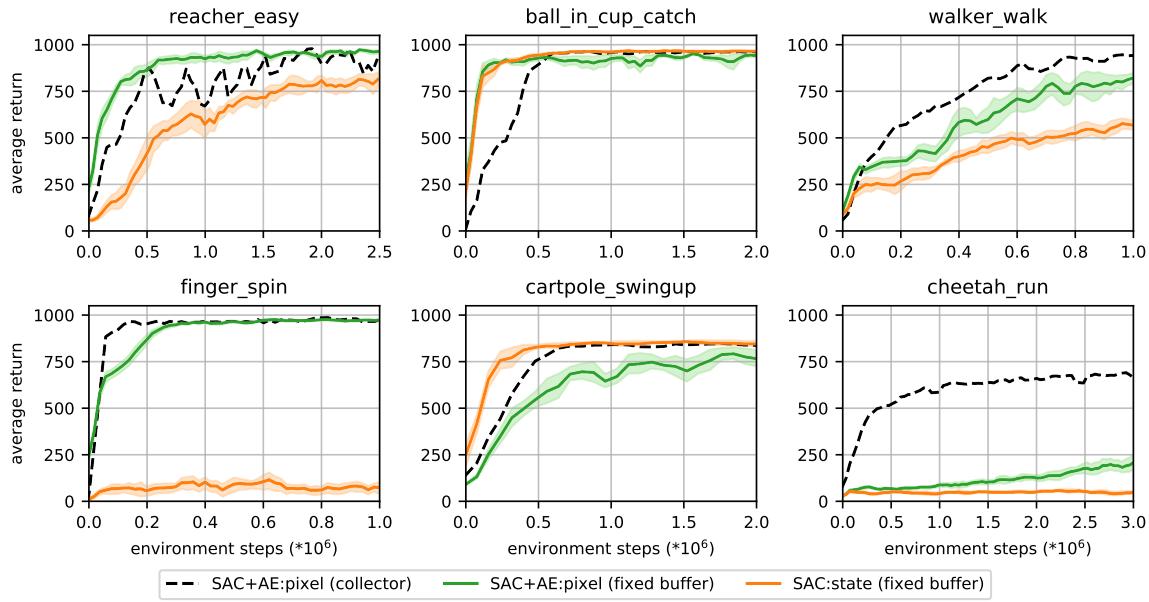


Figure 16. An auxiliary signal is provided by reconstructing a low-dimensional state from the corresponding image observation. Perhaps surprisingly, such *synthetic* supervision doesn't guarantee sufficient signal to fit the high-capacity encoder, which we infer from the suboptimal performance of SAC:pixel (state supervision) compared to SAC:pixel in `ball_in_cup_catch`.

## I. Optimality of Learned Latent Representation

We define the optimality of the learned latent representation as the ability of our model to extract and preserve all relevant information from the pixel observations sufficient to learn a good policy. For example, the proprioceptive state representation is clearly better than the pixel representation because we can learn a better policy. However, the differences in performance of SAC:state and SAC+AE:pixel can be attributed not only to the different observation spaces, but also the difference in data collected in the replay buffer. To decouple these attributes and determine how much information loss there is in moving from proprioceptive state to pixel images, we measure final task reward of policies learned from the same fixed replay buffer, where one is trained on proprioceptive states and the other trained on pixel observations.

We first train a SAC+AE policy until convergence and save the replay buffer that we collected during training. Importantly, in the replay buffer we store both the pixel observations and the corresponding proprioceptive states. Note that for two policies trained on the fixed replay buffer, we are operating in an off-policy regime, and thus it is possible we won't be able to train a policy that performs as well.



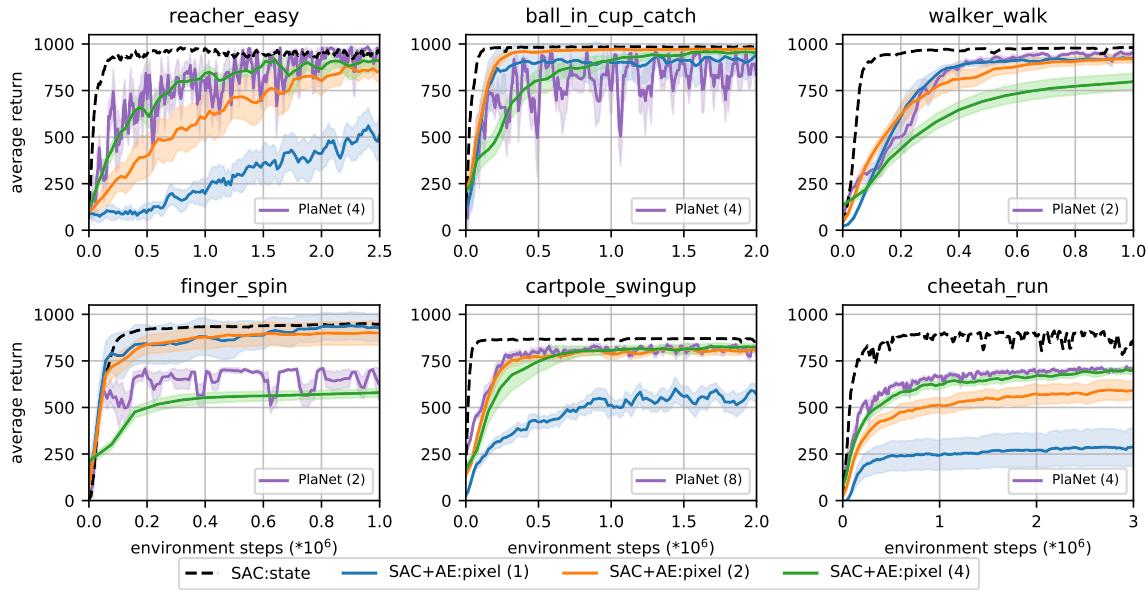
*Figure 17.* Training curves for the policy used to collect the buffer (*SAC+AE:pixel (collector)*), and the two policies learned on that buffer using proprioceptive (*SAC:state (fixed buffer)*) and pixel observations (*SAC+AE:pixel (fixed buffer)*). We see that our method actually outperforms proprioceptive observations in this setting.

In Figure 17 we find, surprisingly, that our learned latent representation outperforms proprioceptive state on a fixed buffer. This could be because the data collected in the buffer is by a policy also learned from pixel observations, and is different enough from the policy that would be learned from proprioceptive states that SAC:state underperforms in this setting.

## J. Importance of Action Repeat

We found that repeating nominal actions several times has a significant effect on learning dynamics and final reward. Prior works (Hafner et al., 2018; Lee et al., 2019) treat action repeat as a hyper parameter to the learning algorithm, rather than a property of the target environment. Effectively, action repeat decreases the control horizon of the task and makes the control dynamics more stable. Yet, action repeat can also introduce a harmful bias, that prevents the agent from learning an optimal policy due to the injected lag. This tasks a practitioner with a problem of finding an optimal value for the action repeat hyper parameter that stabilizes training without limiting control elasticity too much.

To get more insights, we perform an ablation study, where we sweep over several choices for action repeat on multiple control tasks and compare acquired results against PlaNet (Hafner et al., 2018) with the original action repeat setting, which was also tuned per environment. We use the same setup as detailed in Appendix B.3. Specifically, we average performance over 10 random seeds, and reduce the number of training observations inverse proportionally to the action repeat value. The results are shown in Figure 18. We observe that PlaNet’s choice of action repeat is not always optimal for our algorithm. For example, we can significantly improve performance of our agent on the ball\_in\_cup\_catch task if instead of taking the same nominal action four times, as PlaNet suggests, we take it once or twice. The same is true on a few other environments.



*Figure 18.* We study the importance of the action repeat hyper parameter on final performance. We evaluate three different settings, where the agent applies a sampled action once (SAC+AE:pixel (1)), twice (SAC+AE:pixel (2)), or four times (SAC+AE:pixel (4)). As a reference, we also plot the PlaNet (Hafner et al., 2018) results with the original action repeat setting. Action repeat has a significant effect on learning. Moreover, we note that the PlaNet’s choice of hyper parameters is not always optimal for our method (e.g. it is better to apply an action only once on walker\_walk, than taking it twice).