

---

# Stabilizing Deep $Q$ -Learning with ConvNets and Vision Transformers under Data Augmentation

---

Nicklas Hansen<sup>1</sup> Hao Su<sup>1</sup> Xiaolong Wang<sup>1</sup>

<sup>1</sup>University of California, San Diego

nihansen@ucsd.edu {haosu,xiw012}@eng.ucsd.edu

## Abstract

While agents trained by Reinforcement Learning (RL) can solve increasingly challenging tasks directly from visual observations, generalizing learned skills to novel environments remains very challenging. Extensive use of data augmentation is a promising technique for improving generalization in RL, but it is often found to decrease sample efficiency and can even lead to divergence. In this paper, we investigate causes of instability when using data augmentation in common off-policy RL algorithms. We identify two problems, both rooted in high-variance  $Q$ -targets. Based on our findings, we propose a simple yet effective technique for stabilizing this class of algorithms under augmentation. We perform extensive empirical evaluation of image-based RL using both ConvNets and Vision Transformers (ViT) on a family of benchmarks based on DeepMind Control Suite, as well as in robotic manipulation tasks. Our method greatly improves stability and sample efficiency of ConvNets under augmentation, and achieves generalization results competitive with state-of-the-art methods for image-based RL in environments with unseen visuals. We further show that our method scales to RL with ViT-based architectures, and that data augmentation may be especially important in this setting.<sup>†</sup>

## 1 Introduction

Reinforcement Learning (RL) from visual observations has achieved tremendous success in various applications such as video-games [43, 4, 70], robotic manipulation [37], and autonomous navigation [42, 83]. However, it is still very challenging for current methods to generalize the learned skills to novel environments, and policies trained by RL can easily overfit to the training environment [81, 13], especially for high-dimensional observation spaces such as images [8, 58].

Increasing the variability in training data via domain randomization [66, 50] and data augmentation [57, 35, 33, 51] has demonstrated encouraging results for learning policies invariant to changes in environment observations. Specifically, recent works on data augmentation [35, 33] both show improvements in sample efficiency from simple cropping and translation augmentations, but the studies also conclude that additional data augmentation in fact *decrease* sample efficiency and even cause divergence. While these augmentations have the potential to improve generalization, the increasingly varied data makes the optimization more challenging and risks instability. Unlike supervised learning, balancing the trade-off between stability and generalization in RL requires substantial trial and error.

In this paper, we illuminate causes of instability when applying data augmentation to common off-policy RL algorithms [43, 38, 15, 18]. Based on our findings, we provide an intuitive method for stabilizing this class of algorithms under use of strong data augmentation. Specifically, we find two main causes of instability in previous work’s application of data augmentation: (i) indiscriminate application of data augmentation resulting in high-variance  $Q$ -targets; and (ii) that  $Q$ -value estimation strictly from augmented data results in over-regularization.

---

<sup>†</sup>Website and code is available at: <https://nicklashansen.github.io/SVEA>.

To address these problems, we propose **SVEA: Stabilized  $Q$ -Value Estimation under Augmentation**, a simple yet effective framework for data augmentation in off-policy RL that greatly improves stability of  $Q$ -value estimation. Our method consists of the following three components: Firstly, by only applying augmentation in  $Q$ -value estimation of the *current* state, *without* augmenting  $Q$ -targets used for bootstrapping, SVEA circumvents erroneous bootstrapping caused by data augmentation; Secondly, we formulate a modified  $Q$ -objective that optimizes  $Q$ -value estimation jointly over both augmented and unaugmented copies of the observations; Lastly, for SVEA implemented with an actor-critic algorithm, we optimize the actor strictly on unaugmented data, and instead learn a generalizable policy indirectly through parameter-sharing. Our framework can be implemented efficiently without additional forward passes nor introducing additional learnable parameters.

We perform extensive empirical evaluation on the DeepMind Control Suite [64] and extensions of it, including the DMControl Generalization Benchmark [21] and the Distracting Control Suite [60], as well as a set of robotic manipulation tasks. Our method greatly improve  $Q$ -value estimation with ConvNets under a set of strong data augmentations, and achieves sample efficiency, asymptotic performance, and generalization that is competitive or better than previous state-of-the-art methods in all tasks considered, at a lower computational cost. Finally, we show that our method scales to RL with Vision Transformers (ViT) [10]. We find that ViT-based architectures are especially prone to overfitting, and data augmentation may therefore be a key component for large-scale RL.

## 2 Related Work

**Representation Learning.** Learning visual invariances using data augmentation and self-supervised objectives has proven highly successful in computer vision [46, 45, 82, 74, 68, 65, 75, 27, 7]. For example, Chen et al. [7] perform an extensive study on data augmentation (e.g. random cropping and image distortions) for contrastive learning, and show that representations pre-trained with such transformations transfer effectively to downstream tasks. While our work also uses data augmentation for learning visual invariances, we leverage the  $Q$ -objective of deep  $Q$ -learning algorithms instead of auxiliary representation learning tasks.

**Visual Learning for RL.** Numerous methods have been proposed with the goal of improving sample efficiency [29, 56, 68, 76, 40, 59, 61, 54, 77] of image-based RL. Recently, using self-supervision to improve generalization in RL has also gained interest [80, 47, 55, 1, 22, 21, 72]. Notably, Zhang et al. [80] and Agarwal et al. [1] propose to learn behavioral similarity embeddings via auxiliary tasks (bisimulation metrics and contrastive learning, respectively), and Hansen et al. [21] learn visual invariances through an auxiliary prediction task. While these results are encouraging, it has also been shown in [29, 40, 22, 79, 41] that the best choice of auxiliary tasks depends on the particular RL task, and that joint optimization with sub-optimally chosen tasks can lead to gradient interference. We achieve competitive sample-efficiency and generalization results *without* the need for carefully chosen auxiliary tasks, and our method is therefore applicable to a larger variety of RL tasks.

**Data Augmentation and Randomization for RL.** Our work is directly inspired by previous work on generalization in RL by domain randomization [66, 50, 48, 52, 6] and data augmentation [36, 9, 71, 35, 33, 51, 61, 21]. For example, Tobin et al. [66] show that a neural network trained for object localization in a simulation with randomized visual augmentations improves real world generalization. Similarly, Lee et al. [36] show that application of a random convolutional layer to observations during training improve generalization in 3D navigation tasks. More recently, extensive studies on data augmentation [35, 33] have been conducted with RL, and conclude that, while small random crops and translations can improve sample efficiency, most data augmentations *decrease* sample efficiency and cause divergence. We illuminate main causes of instability, and propose a framework for data augmentation in deep  $Q$ -learning algorithms that drastically improves stability and generalization.

**Improving Deep  $Q$ -Learning.** While deep  $Q$ -learning algorithms such as Deep  $Q$ -Networks (DQN) [43] have achieved impressive results in image-based RL, the temporal difference objective is known to have inherent instabilities when used in conjunction with function approximation and off-policy data [63]. Therefore, a variety of algorithmic improvements have been proposed to improve convergence [24, 73, 25, 23, 53, 38, 15, 14, 28]. For example, Hasselt et al. [24] reduce overestimation of  $Q$ -values by decomposing the target  $Q$ -value estimation into action selection and action evaluation using separate networks. Lillicrap et al. [38] reduce target variance by defining the target  $Q$ -network as a slow-moving average of the online  $Q$ -network. Our method also improves  $Q$ -value estimation, but we specifically address the instability of deep  $Q$ -learning algorithms on augmented data.

### 3 Preliminaries

**Problem formulation.** We formulate the interaction between environment and policy as a Markov Decision Process (MDP) [2]  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\mathcal{P}: \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  is the state transition function that defines a conditional probability distribution  $\mathcal{P}(\cdot | \mathbf{s}_t, \mathbf{a}_t)$  over all possible next states given a state  $\mathbf{s}_t \in \mathcal{S}$  and action  $\mathbf{a}_t \in \mathcal{A}$  taken at time  $t$ ,  $r: \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$  is a reward function, and  $\gamma \in [0, 1)$  is the discount factor. Because image observations only offer partial state observability [30], we define a state  $\mathbf{s}_t$  as a sequence of  $k + 1$  consecutive frames  $(\mathbf{o}_t, \mathbf{o}_{t-1}, \dots, \mathbf{o}_{t-k})$ ,  $\mathbf{o} \in \mathcal{O}$ , where  $\mathcal{O}$  is the high-dimensional image space, as proposed in Mnih et al. [43]. The goal is then to learn a policy  $\pi: \mathcal{S} \mapsto \mathcal{A}$  that maximizes discounted return  $R_t = \mathbb{E}_{\Gamma \sim \pi} [\sum_{t=1}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$  along a trajectory  $\Gamma = (\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_T)$  obtained by following policy  $\pi$  from an initial state  $\mathbf{s}_0 \in \mathcal{S}$  to a state  $\mathbf{s}_T$  with state transitions sampled from  $\mathcal{P}$ , and  $\pi$  is parameterized by a collection of learnable parameters  $\theta$ . For clarity, we hereon generically denote parameterization with subscript, e.g.  $\pi_\theta$ . We further aim to learn parameters  $\theta$  s.t.  $\pi_\theta$  generalizes well (i.e., obtains high discounted return) to unseen MDPs, which is generally unfeasible without further assumptions about the structure of the space of MDPs. In this work, we focus on generalization to MDPs  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$ , where states  $\overline{\mathbf{s}}_t \in \overline{\mathcal{S}}$  are constructed from observations  $\overline{\mathbf{o}}_t \in \overline{\mathcal{O}}$ ,  $\mathcal{O} \subseteq \overline{\mathcal{O}}$  of a *perturbed* observation space  $\overline{\mathcal{O}}$  (e.g. unseen visuals), and  $\overline{\mathcal{M}} \sim \mathbb{M}$  for a space of MDPs  $\mathbb{M}$ .

**Deep Q-Learning.** Common model-free off-policy RL algorithms aim to estimate an optimal state-action value function  $Q^*: \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$  as  $Q_\theta(\mathbf{s}, \mathbf{a}) \approx Q^*(\mathbf{s}, \mathbf{a}) = \max_{\pi_\theta} \mathbb{E}[R_t | \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}]$  using function approximation. In practice, this is achieved by means of the single-step Bellman residual  $\left( r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}'_t} Q_\psi^{\text{tgt}}(\mathbf{s}_{t+1}, \mathbf{a}'_t) \right) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t)$  [62], where  $\psi$  parameterizes a *target* state-action value function  $Q^{\text{tgt}}$ . We can choose to minimize this residual (also known as the *temporal difference error*) directly wrt  $\theta$  using a mean squared error loss, which gives us the objective

$$\mathcal{L}_Q(\theta, \psi) = \mathbb{E}_{\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1} \sim \mathcal{B}} \left[ \frac{1}{2} \left[ \left( r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}'_t} Q_\psi^{\text{tgt}}(\mathbf{s}_{t+1}, \mathbf{a}'_t) \right) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t) \right]^2 \right], \quad (1)$$

where  $\mathcal{B}$  is a replay buffer with transitions collected by a behavioral policy [39]. From here, we can derive a greedy policy directly by selecting actions  $\mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\theta(\mathbf{s}_t, \mathbf{a}_t)$ . While  $Q^{\text{tgt}} = Q$  and periodically setting  $\psi \leftarrow \theta$  exactly recovers the objective of DQN [43], several improvements have been proposed to improve stability of Eq. 1, such as Double Q-learning [24], Dueling Q-networks [73], updating target parameters using a slow-moving average of the online Q-network [38]:

$$\psi_{n+1} \leftarrow (1 - \zeta)\psi_n + \zeta\theta_n \quad (2)$$

for an iteration step  $n$  and a momentum coefficient  $\zeta \in (0, 1)$ , and others [25, 23, 53, 14, 28]. As computing  $\max_{\mathbf{a}'_t} Q_\psi^{\text{tgt}}(\mathbf{s}_{t+1}, \mathbf{a}'_t)$  in Eq. 1 is intractable for large and continuous action spaces, a number of prominent *actor-critic* algorithms that additionally learn a policy  $\pi_\theta(\mathbf{s}_t) \approx \arg \max_{\mathbf{a}_t} Q_\theta(\mathbf{s}_t, \mathbf{a}_t)$  have therefore been proposed [38, 15, 18].

**Soft Actor-Critic (SAC)** [18] is an off-policy actor-critic algorithm that learns a state-action value function  $Q_\theta$  and a stochastic policy  $\pi_\theta$  (and optionally a temperature parameter), where  $Q_\theta$  is optimized using a variant of the objective in Eq. 1 and  $\pi_\theta$  is optimized using a  $\gamma$ -discounted maximum-entropy objective [84]. To improve stability, SAC is also commonly implemented using Double Q-learning and the slow-moving target parameters from Eq. 2. We will in the remainder of this work describe our method in the context of a generic off-policy RL algorithm that learns a parameterized state-action value function  $Q_\theta$ , while we in our experiments discussed in Section 6 evaluate of our method using SAC as base algorithm.

### 4 Pitfalls of Data Augmentation in Deep Q-Learning

In this section, we aim to illuminate the main causes of instability from naïve application of data augmentation in Q-value estimation. Our goal is to learn a Q-function  $Q_\theta$  for an MDP  $\mathcal{M}$  that generalizes to novel MDPs  $\overline{\mathcal{M}} \sim \mathbb{M}$  with unseen visuals, and we leverage data augmentation as an optimality-invariant state transformation  $\tau$  to induce a bisimulation relation [34, 17] between a state  $\mathbf{s}$  and its transformed (augmented) counterpart  $\mathbf{s}^{\text{aug}} = \tau(\mathbf{s}, \nu)$  with parameters  $\nu \sim \mathcal{V}$ .

**Definition 1** (Optimality-Invariant State Transformation [33]). *Given an MDP  $\mathcal{M}$ , a state transformation  $\tau: \mathcal{S} \times \mathcal{V} \mapsto \mathcal{S}$  is an optimality-invariant state transformation if  $Q(\mathbf{s}, \mathbf{a}) = Q(\tau(\mathbf{s}, \nu), \mathbf{a}) \quad \forall \mathbf{s} \in \mathcal{S}, \mathbf{a} \in \mathcal{A}, \nu \in \mathcal{V}$ , where  $\nu \in \mathcal{V}$  parameterizes the transformation  $\tau$ .*

Following our definitions of  $\mathcal{M}, \overline{\mathcal{M}}$  from Section 3, we can further extend the concept of optimality-invariant transformations to MDPs, noting that a change of state space (e.g. perturbed visuals) itself can be described as a transformation  $\overline{\tau}: \mathcal{S} \times \overline{\mathcal{V}} \mapsto \overline{\mathcal{S}}$  with unknown parameters  $\overline{\nu} \in \overline{\mathcal{V}}$ . If we choose the set of parameters  $\mathcal{V}$  of a state transformation  $\tau$  to be sufficiently large such that it intersects with  $\overline{\mathcal{V}}$  with high probability, we can therefore expect to improve generalization to state and observation spaces not seen during training. However, while naïve application of data augmentation as in previous work [35, 33, 61, 54] may potentially improve generalization, it can be harmful to  $Q$ -value estimation. We hypothesize that this is primarily because it dramatically increases the size of the observed state space, and consequently also increases variance  $\text{Var}[Q(\tau(\mathbf{s}, \nu))] \geq \text{Var}[Q(\mathbf{s})]$ ,  $\nu \sim \mathcal{V}$  when  $\mathcal{V}$  is large. Concretely, we identify the following two issues:

**Pitfall 1: Non-deterministic  $Q$ -target.** For deep  $Q$ -learning algorithms, previous work [35, 33, 61, 54] applies augmentation to both state  $\mathbf{s}_t^{\text{aug}} \triangleq \tau(\mathbf{s}_t, \nu)$  and successor state  $\mathbf{s}_{t+1}^{\text{aug}} \triangleq \tau(\mathbf{s}_{t+1}, \nu')$  where  $\nu, \nu' \sim \mathcal{V}$ . Compared with DQN [43] that uses a deterministic (more precisely, periodically updated)  $Q$ -target, this practice introduces a non-deterministic  $Q$ -target  $r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}'_t} Q_{\psi}^{\text{tgt}}(\mathbf{s}_{t+1}^{\text{aug}}, \mathbf{a}'_t)$  depending on the augmentation parameters  $\nu'$ . As observed in the original DQN paper, high-variance target values are detrimental to  $Q$ -learning algorithms, and may cause divergence due to the “deadly triad” of function approximation, bootstrapping, and off-policy learning [63]. This motivates the work to introduce a slowly changing target network, and several other works have refined the  $Q$ -target update rule [38, 15] to further reduce volatility. However, because data augmentation is inherently non-deterministic, it can greatly increase variance in  $Q$ -target estimation and exacerbates the issue of volatility, as shown in Figure 1 (top). This is particularly troubling in actor-critic algorithms such as DDPG [38] and SAC [18], where the  $Q$ -target is estimated from  $(\mathbf{s}_{t+1}, \mathbf{a}')$ ,  $\mathbf{a}' \sim \pi(\cdot | \mathbf{s}_{t+1})$ , which introduces an additional source of error from  $\pi$  that is non-negligible especially when  $\mathbf{s}_{t+1}$  is augmented.

**Pitfall 2: Over-regularization.** Data augmentation was originally introduced in the supervised learning regime as a regularizer to prevent overfitting of high-capacity models. However, for RL, even learning a policy in the training environment is hard. While data augmentation may improve generalization, it greatly increases the difficulty of policy learning, i.e., optimizing  $\theta$  for  $Q_{\theta}$  and potentially a behavior network  $\pi_{\theta}$ . Particularly, when the temporal difference loss from Eq. 1 cannot be well minimized, the large amount of augmented states dominate the gradient, which significantly impacts  $Q$ -value estimation of both augmented and unaugmented states. We refer to this issue as *over-regularization* by data augmentation. Figure 1 (bottom) shows the mean difference in  $Q$ -predictions made with augmented vs. unaugmented data in fully converged DrQ [33] agents trained with *shift* augmentation. Augmentations such as affine-jitter, random convolution, and random overlay incur large differences in estimated  $Q$ -values. While such difference can be reduced by regularizing the optimization with each individual augmentation, we emphasize that even the minimal shift augmentation used throughout training incurs non-zero difference. Since  $\psi$  is commonly chosen to be a moving average of  $\theta$  as in Eq. 2, such differences caused by over-regularization affect  $Q_{\theta}$  and  $Q_{\psi}^{\text{tgt}}$  equally, and optimization may therefore still diverge depending on the choice of data augmentation. As such, there is an inherent trade-off between accurate  $Q$ -value estimation and generalization when using data augmentation. In the following section, we address these pitfalls.

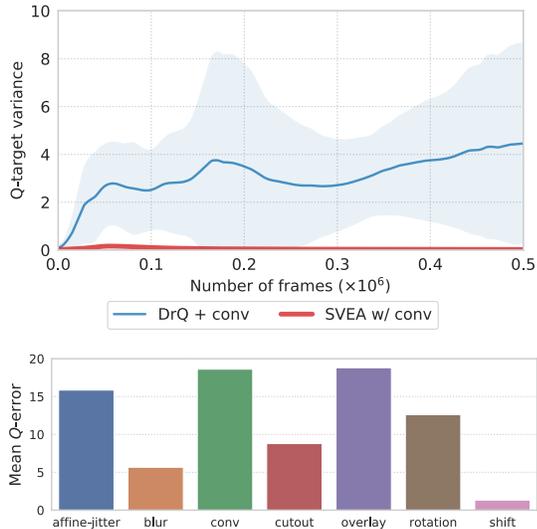


Figure 1. (Top) Mean  $Q$ -target variance of DrQ [33] and SVEA (ours), both trained with *conv* augmentation [36]. (Bottom) Mean difference in  $Q$ -value estimation on augmented vs. non-augmented data. We measure mean absolute error in  $Q$ -value estimation from converged DrQ agents (trained with *shift* augmentation) on the same observations before and after augmentation. Both figures are averages across 5 seeds for each of the 5 tasks from DMControl-GB.

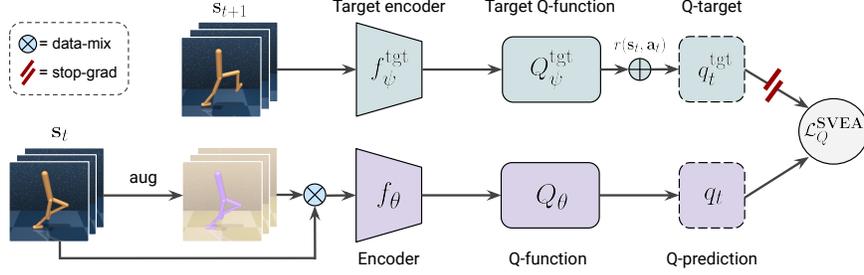


Figure 2. **Overview.** An observation  $s_t$  is transformed by data augmentation  $\tau(\cdot, \nu)$ ,  $\nu \sim \mathcal{V}$  to produce a view  $s_t^{\text{aug}}$ . The  $Q$ -function  $Q_{\theta}$  is then jointly optimized on both augmented and unaugmented data wrt the objective in Eq. 7, with the  $Q$ -target of the Bellman equation computed from an unaugmented observation  $s_{t+1}$ . We illustrate our data-mixing strategy by the  $\otimes$  operator.

## 5 Method

We propose **SVEA: Stabilized  $Q$ -Value Estimation under Augmentation**, a general framework for visual generalization in RL by use of data augmentation. SVEA applies data augmentation in a novel learning framework leveraging two data streams – with and without augmented data, respectively. Our method is compatible with any standard off-policy RL algorithm without changes to the underlying neural network that parameterizes the policy, and it requires no additional forward passes, auxiliary tasks, nor learnable parameters. While SVEA in principle does not make any assumptions about the structure of states  $s_t \in \mathcal{S}$ , we here describe our method in the context of image-based RL.

### 5.1 Architectural Overview

An overview of the SVEA architecture is provided in Figure 2. Our method leverages properties of common neural network architectures used in off-policy RL without introducing additional learnable parameters. We subdivide the neural network layers and corresponding learnable parameters of a state-action value function into sub-networks  $f_{\theta}$  (denoted the state *encoder*) and  $Q_{\theta}$  (denoted the  *$Q$ -function*) s.t.  $q_t \triangleq Q_{\theta}(f_{\theta}(s_t), \mathbf{a}_t)$  is the predicted  $Q$ -value corresponding to a given state-action pair  $(s_t, \mathbf{a}_t)$ . We similarly define the target state-action value function s.t.  $q_t^{\text{tgt}} \triangleq r(s_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}'} Q_{\psi}^{\text{tgt}}(f_{\psi}^{\text{tgt}}(s_{t+1}), \mathbf{a}')$  is the target  $Q$ -value for  $(s_t, \mathbf{a}_t)$ , and we define parameters  $\psi$  as an exponential moving average of  $\theta$  as in Eq. 2. Depending on the choice of underlying algorithm, we may choose to additionally learn a parameterized policy  $\pi_{\theta}$  that shares encoder parameters with  $Q_{\theta}$  and selects actions  $\mathbf{a}_t \sim \pi_{\theta}(\cdot | f_{\theta}(s_t))$ .

To circumvent erroneous bootstrapping from augmented data (as discussed in Section 4), we strictly apply data augmentation in  $Q$ -value estimation of the *current* state  $s_t$ , *without* applying data augmentation to the successor state  $s_{t+1}$  used in Eq. 1 for bootstrapping with  $Q_{\psi}^{\text{tgt}}$  (and  $\pi_{\theta}$  if applicable), which addresses Pitfall 1. If  $\pi_{\theta}$  is learned (i.e., SVEA is implemented with an actor-critic algorithm), we also optimize it strictly from unaugmented data. To mitigate over-regularization in optimization of  $f_{\theta}$  and  $Q_{\theta}$  (Pitfall 2), we further employ a modified  $Q$ -objective that leverages both augmented and unaugmented data, which we introduce in the following section.

### 5.2 Learning Objective

Our method redefines the temporal difference objective from Eq. 1 to better leverage data augmentation. First, recall that  $q_t^{\text{tgt}} = r(s_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}'} Q_{\psi}^{\text{tgt}}(f_{\psi}^{\text{tgt}}(s_{t+1}), \mathbf{a}')$ . Instead of learning to predict  $q_t^{\text{tgt}}$  only from state  $s_t$ , we propose to minimize a nonnegative linear combination of  $\mathcal{L}_Q$  over two individual data streams,  $s_t$  and  $s_t^{\text{aug}} = \tau(s_t, \nu)$ ,  $\nu \sim \mathcal{V}$ , which we define as the objective

$$\mathcal{L}_Q^{\text{SVEA}}(\theta, \psi) \triangleq \alpha \mathcal{L}_Q(s_t, q_t^{\text{tgt}}; \theta, \psi) + \beta \mathcal{L}_Q(s_t^{\text{aug}}, q_t^{\text{tgt}}; \theta, \psi) \quad (3)$$

$$= \mathbb{E}_{s_t, \mathbf{a}_t, s_{t+1} \sim \mathcal{B}} \left[ \alpha \|Q_{\theta}(f_{\theta}(s_t), \mathbf{a}_t) - q_t^{\text{tgt}}\|_2^2 + \beta \|Q_{\theta}(f_{\theta}(s_t^{\text{aug}}), \mathbf{a}_t) - q_t^{\text{tgt}}\|_2^2 \right], \quad (4)$$

where  $\alpha, \beta$  are constant coefficients that balance the ratio of the **unaugmented** and **augmented** data streams, respectively, and  $q_t^{\text{tgt}}$  is computed strictly from unaugmented data.  $\mathcal{L}_Q^{\text{SVEA}}(\theta, \psi)$  serves as a *data-mixing* strategy that oversamples unaugmented data as an implicit variance reduction technique.

As we will verify empirically in Section 6, data-mixing is a simple and effective technique for variance reduction that works well in tandem with our proposed modifications to bootstrapping. For  $\alpha = \beta$ , the objective in Eq. 4 can be evaluated in a single, batched forward-pass by rewriting it as:

$$\mathbf{g}_t = [\mathbf{s}_t, \tau(\mathbf{s}_t, \nu)]_N \quad (5)$$

$$h_t = [q_t^{\text{tgt}}, q_t^{\text{tgt}}]_N \quad (6)$$

$$\mathcal{L}_Q^{\text{SVEA}}(\theta, \psi) = \mathbb{E}_{\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1} \sim \mathcal{B}, \nu \sim \mathcal{V}} \left[ (\alpha + \beta) \|Q_\theta(f_\theta(\mathbf{g}_t), \mathbf{a}_t) - h_t\|_2^2 \right], \quad (7)$$

where  $[\cdot]_N$  is a concatenation operator along the batch dimension  $N$  for  $\mathbf{s}_t, \mathbf{s}_t^{\text{aug}} \in \mathbb{R}^{N \times C \times H \times W}$  and  $q_t^{\text{tgt}} \in \mathbb{R}^{N \times 1}$ , which is illustrated as  $\otimes$  in Figure 2. Empirically, we find  $\alpha = 0.5, \beta = 0.5$  to be both effective and practical to implement, which we adopt in the majority of our experiments. However, more sophisticated schemes for selecting  $\alpha, \beta$  and/or varying them as training progresses could be interesting directions for future research. If the base algorithm of choice learns a policy  $\pi_\theta$ , its objective  $\mathcal{L}_\pi(\theta)$  is optimized solely on unaugmented states  $\mathbf{s}_t$  without changes to the objective, and a stop-grad operation is applied after  $f_\theta$  to prevent non-stationary gradients of  $\mathcal{L}_\pi(\theta)$  from interfering with  $Q$ -value estimation, i.e., only the objective from Eq. 4 or optionally Eq. 7 updates  $f_\theta$  using stochastic gradient descent. As described in Section 5.1, parameters  $\psi$  are updated using an exponential moving average of  $\theta$  and a stop-grad operation is therefore similarly applied after  $Q_\psi^{\text{tgt}}$ . We summarize our method for  $\alpha = \beta$  applied to a generic off-policy algorithm in Algorithm 1.

**Algorithm 1** Generic SVEA off-policy algorithm (► naïve augmentation, ► our modifications)

$\theta, \theta_\pi, \psi$ : randomly initialized network parameters,  $\psi \leftarrow \theta$     ► Initialize  $\psi$  to be equal to  $\theta$   
 $\eta, \zeta$ : learning rate and momentum coefficient  
 $\alpha, \beta$ : loss coefficients, *default*: ( $\alpha = 0.5, \beta = 0.5$ )

1: **for** timestep  $t = 1 \dots T$  **do**  
**act:**  
2:     $\mathbf{a}_t \sim \pi_\theta(\cdot | f_\theta(\mathbf{s}_t))$     ► Sample action from policy  
3:     $\mathbf{s}'_t \sim \mathcal{P}(\cdot | \mathbf{s}_t, \mathbf{a}_t)$     ► Sample transition from environment  
4:     $\mathcal{B} \leftarrow \mathcal{B} \cup (\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}'_t)$     ► Add transition to replay buffer

**update:**  
5:     $\{\mathbf{s}_i, \mathbf{a}_i, r(\mathbf{s}_i, \mathbf{a}_i), \mathbf{s}'_i \mid i = 1 \dots N\} \sim \mathcal{B}$     ► Sample batch of transitions  
6:     $\mathbf{s}_i = \tau(\mathbf{s}_i, \nu_i), \mathbf{s}'_i = \tau(\mathbf{s}'_i, \nu'_i), \nu_i, \nu'_i \sim \mathcal{V}$     ► Naïve application of data augmentation  
7:    **for** transition  $i = 1 \dots N$  **do**  
8:      $\theta_\pi \leftarrow \theta_\pi - \eta \nabla_{\theta_\pi} \mathcal{L}_\pi(\mathbf{s}_i; \theta_\pi)$  (if applicable)    ► Optimize  $\pi_\theta$  with SGD  
9:      $q_i^{\text{tgt}} = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_\psi^{\text{tgt}}(f_\psi^{\text{tgt}}(\mathbf{s}'_i), \mathbf{a}'_i)$     ► Compute  $Q$ -target  
10:      $\mathbf{s}_i^{\text{aug}} = \tau(\mathbf{s}_i, \nu_i), \nu_i \sim \mathcal{V}$     ► Apply stochastic data augmentation  
11:      $\mathbf{g}_i = [\mathbf{s}_i, \mathbf{s}_i^{\text{aug}}]_N, h_i = [q_i^{\text{tgt}}, q_i^{\text{tgt}}]_N$     ► Pack data streams  
12:      $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_Q^{\text{SVEA}}(\mathbf{g}_i, h_i; \theta, \psi)$     ► Optimize  $f_\theta$  and  $Q_\theta$  with SGD  
13:      $\psi \leftarrow (1 - \zeta)\psi + \zeta\theta$     ► Update  $\psi$  using EMA of  $\theta$

## 6 Experiments

We evaluate both sample efficiency, asymptotic performance, and generalization of our method and a set of strong baselines using both ConvNets and Vision Transformers (ViT) [10] in tasks from DeepMind Control Suite (DMControl) [64] as well as a set of robotic manipulation tasks. DMControl offers challenging and diverse continuous control tasks and is widely used as a benchmark for image-based

RL [19, 20, 76, 59, 35, 33]. To evaluate generalization of our method and baselines, we test methods under challenging distribution shifts (as illustrated in Figure 3) from the DMControl Generalization Benchmark (DMControl-GB) [21], the Distracting Control Suite (DistractingCS) [60], as well as distribution shifts unique to the robotic manipulation environment. Code is available at <https://github.com/nicklashansen/dmcontrol-generalization-benchmark>.

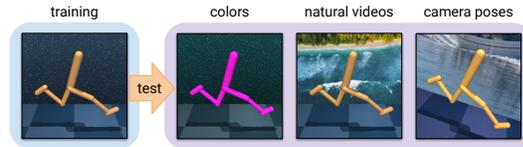
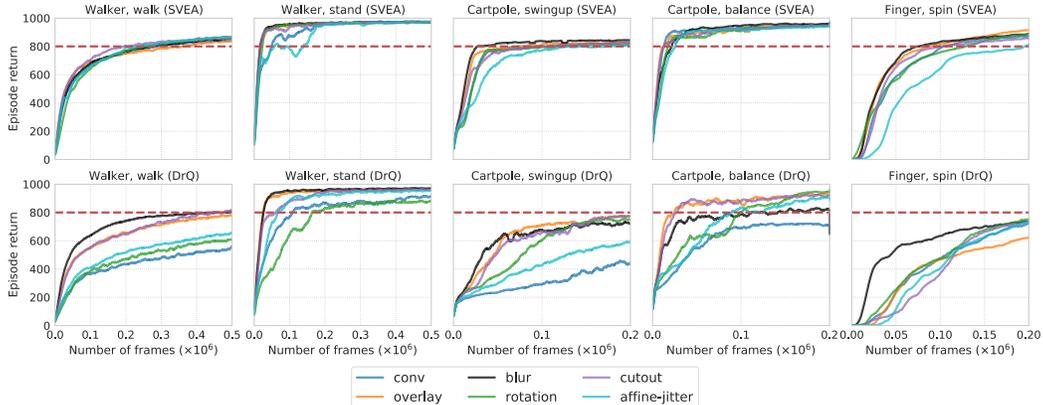
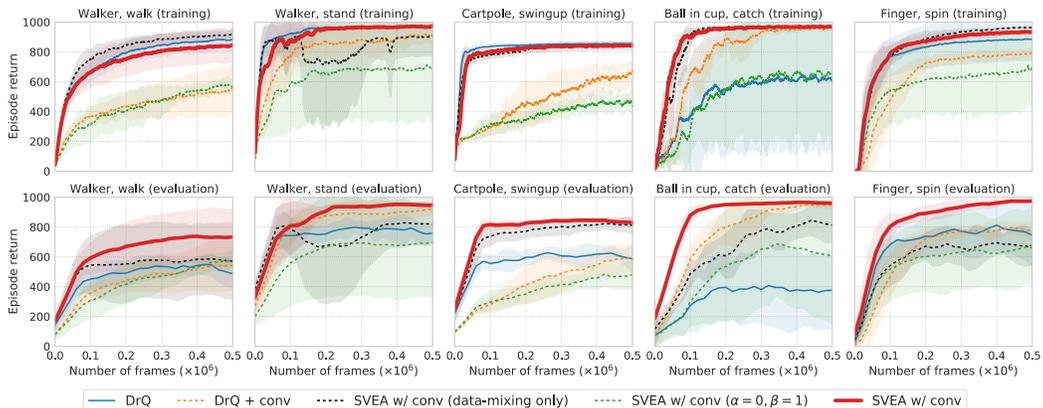


Figure 3. **Experimental setup.** Agents are trained in a fixed environment and are expected to generalize to novel environments with e.g. random colors, backgrounds, and camera poses.



**Figure 4. Data augmentations.** Training performance of SVEA (top) and DrQ (bottom) under 6 common data augmentations. Mean of 5 seeds. Red line at 800 return is for visual guidance only. We omit visualization of std. deviations for clarity, but provide per-augmentation comparisons to DrQ (including std. deviations) across all tasks in Appendix B, and test performances in Appendix C.



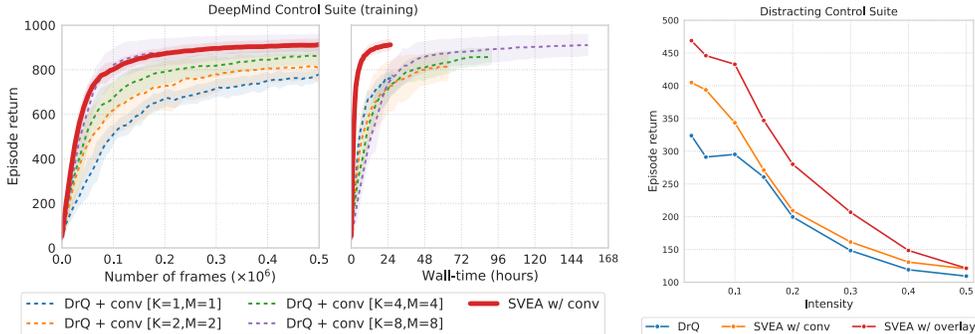
**Figure 5. Training and test performance.** We compare SVEA to DrQ with and without random convolution augmentation, as well as a set of ablations. *Data-mixing only* indiscriminately applies our data-mixing strategy to all data streams, and  $(\alpha = 0, \beta = 1)$  only augments  $Q$ -predictions but without data-mixing. We find both components to contribute to SVEA’s success. *Top*: episode return on the training environment during training. *Bottom*: generalization measured by episode return on the `color_hard` benchmark of DMControl-GB. Mean of 5 seeds, shaded area is  $\pm 1$  std. deviation.

**Setup.** We implement our method and baselines using SAC [18] as base algorithm, and we apply random shift augmentation to all methods by default. This makes our base algorithm equivalent to DrQ [33] when  $K=1, M=1$ ; we refer to the base algorithm as *unaugmented* and consider stability under additional data augmentation. We use the **same** network architecture and hyperparameters for **all** methods (whenever applicable), and adopt the setup from Hansen and Wang [21]. Observations are stacks of 3 RGB frames of size  $84 \times 84 \times 3$  (and  $96 \times 96 \times 3$  in ViT experiments). In the DMControl-GB and DistractingCS benchmarks, all methods are trained for 500k frames and evaluated on all 5 tasks from DMControl-GB used in prior work, and we adopt the same experimental setup for robotic manipulation. See Appendix H for hyperparameters and further details on our experimental setup.

**Baselines and data augmentations.** We benchmark our method against the following strong baselines: (1) **CURL** [59], a contrastive learning method for RL; (2) **RAD** that applies a random crop; (3) **DrQ** that applies a random shift; (4) **PAD** [22] that adapts to test environments using self-supervision; (5) **SODA** [21] that applies data augmentation in auxiliary learning; as well as a number of ablations. We compare to the  $K=1, M=1$  setting of DrQ by default, but also provide comparison to varying  $K, M$ . We experiment with a diverse set of data augmentations proposed in previous work on RL and computer vision, namely random *shift* [33], random convolution (denoted *conv*) [36], random *overlay* [21], random *cutout* [9], Gaussian *blur*, random *affine-jitter*, and random *rotation* [35, 16]. We provide samples for all data augmentations in Appendix C and test environments in Appendix E.

**Table 1. Comparison to state-of-the-art.** Test performance (episode return) of methods trained in a single, fixed environment and evaluated on (i) randomized colors, and (ii) natural video backgrounds from DMControl-GB. Results for CURL, RAD, PAD, and SODA are obtained from [21] and we report mean and std. deviation over 5 seeds. DrQ corresponds to our SAC base algorithm using random shift augmentation. SVEA matches or outperforms prior methods in all tasks considered.

DMControl-GB (random colors)	CURL	RAD	DrQ	PAD	SODA (conv)	SODA (overlay)	SVEA (conv)	SVEA (overlay)
walker, walk	445 ±99	400 ±61	520 ±91	468 ±47	697 ±66	692 ±68	<b>760</b> ±145	749 ±61
walker, stand	662 ±54	644 ±88	770 ±71	797 ±46	930 ±12	893 ±12	<b>942</b> ±26	933 ±24
cartpole, swingup	454 ±110	590 ±53	586 ±52	630 ±63	831 ±21	805 ±28	<b>837</b> ±23	832 ±23
ball_in_cup, catch	231 ±92	541 ±29	365 ±210	563 ±50	892 ±37	949 ±19	<b>961</b> ±7	959 ±5
finger, spin	691 ±12	667 ±154	776 ±134	803 ±72	901 ±51	793 ±128	<b>977</b> ±5	972 ±6
DMControl-GB (natural videos)	CURL	RAD	DrQ	PAD	SODA (conv)	SODA (overlay)	SVEA (conv)	SVEA (overlay)
walker, walk	556 ±133	606 ±63	682 ±89	717 ±79	635 ±48	768 ±38	612 ±144	<b>819</b> ±71
walker, stand	852 ±75	745 ±146	873 ±83	935 ±20	903 ±56	955 ±13	795 ±70	<b>961</b> ±8
cartpole, swingup	404 ±67	373 ±72	485 ±105	521 ±76	474 ±143	758 ±62	606 ±85	<b>782</b> ±27
ball_in_cup, catch	316 ±119	481 ±26	318 ±157	436 ±55	539 ±111	<b>875</b> ±56	659 ±110	871 ±106
finger, spin	502 ±19	400 ±64	533 ±119	691 ±80	363 ±185	695 ±97	764 ±86	<b>808</b> ±33



**Figure 6. (Left) Comparison with additional DrQ baselines.** We compare SVEA implemented with DrQ [K=1, M=1] as base algorithm to DrQ with varying values of its  $K, M$  hyperparameters. All methods use the *conv* augmentation (in addition to *shift* augmentation used by DrQ). Results are averaged over 5 seeds for each of the 5 tasks from DMControl-GB [21] and shaded area is  $\pm 1$  std. deviation across seeds. Increasing values of  $K, M$  improve sample efficiency of DrQ, but at a high computational cost; DrQ uses approx. 6x wall-time to match the sample efficiency of SVEA. **(Right) DistractingCS.** Episode return as a function of randomization intensity at test-time, aggregated across 5 seeds for each of the 5 tasks from DMControl-GB. See Appendix E for per-task comparison.

### 6.1 Stability and Generalization on DMControl

**Stability.** We evaluate the stability of SVEA and DrQ under 6 common data augmentations; results are shown in Figure 4. While the sample efficiency of DrQ degrades substantially for most augmentations, SVEA is relatively unaffected by the choice of data augmentation and improves sample efficiency in 27 out of 30 instances. While the sample efficiency of DrQ can be improved by increasing its  $K, M$  parameters, we find that DrQ requires approx. 6x wall-time to match the sample efficiency of SVEA; see Figure 6 (left). We further ablate each component of SVEA and report both training and test curves in Figure 5; we find that both components are key to SVEA’s success. Because we empirically find the *conv* augmentation to be particularly difficult to optimize, we provide additional stability experiments in Section 6.2 and 6.3 using this augmentation. See Appendix A for additional ablations.

**Generalization.** We compare the test performance of SVEA to 5 recent state-of-the-art methods for image-based RL on the `color_hard` and `video_easy` benchmarks from DMControl-GB (results in Table 1), as well as the extremely challenging DistractingCS benchmark, where camera pose, background, and colors are continually changing throughout an episode (results in Figure 6 (right)). We here use *conv* and *overlay* augmentations for fair comparison to SODA, and we report additional results on the `video_hard` benchmark in Appendix F. SVEA outperforms all methods considered in 12 out of 15 instances on DMControl-GB, and at a lower computational cost than CURL, PAD, and SODA that all learn auxiliary tasks. On DistractingCS, we observe that SVEA improves generalization by 42% at low intensity, and its generalization degrades significantly slower than DrQ for high intensities. While generalization depends on the particular choice of data augmentation and test environments, this is an encouraging result considering that SVEA enables efficient policy learning with stronger augmentations than previous methods.

## 6.2 RL with Vision Transformers

Vision Transformers (ViT) [10] have recently achieved impressive results on downstream tasks in computer vision. We replace all convolutional layers from the previous experiments with a 4-layer ViT encoder that operates on raw pixels in  $8 \times 8$  space-time patches, and evaluate our method using data augmentation in conjunction with ViT encoders. Importantly, we design the ViT architecture such that it roughly matches our CNN encoder in terms of learnable parameters. The ViT encoder is trained from scratch using RL, and we use the same experimental setup as in our ConvNet experiments. In particular, it is worth emphasizing that both our ViT and CNN encoders are trained using Adam [32] as optimizer and without weight decay. See Figure 7 (top) for an architectural overview, and refer to Appendix H for additional implementation details.

Our training and test results are shown in Figure 7 (bottom). We are, to the best of our knowledge, the first to successfully solve image-based RL tasks without CNNs. We observe that DrQ overfits significantly to the training environment compared to its CNN counterpart (94 test return on `color_hard` for DrQ with ViT vs. 569 with a ConvNet on the *Walker, walk* task). SVEA achieves comparable sample efficiency and improves generalization by 706% and 233% on *Walker, walk* and *Cartpole, swingup*, respectively, over DrQ, while DrQ + conv remains unstable. Interestingly, we observe that our ViT-based implementation of SVEA achieves a mean episode return of 877 on the `color_hard` benchmark of the challenging *Walker, walk* task (vs. 760 using CNNs). SVEA might therefore be a promising technique for future research on RL with CNN-free architectures, where data augmentation appears to be especially important for generalization. We provide additional experiments with ViT encoders in Section 6.3 and make further comparison to ConvNet encoders in Appendix A.

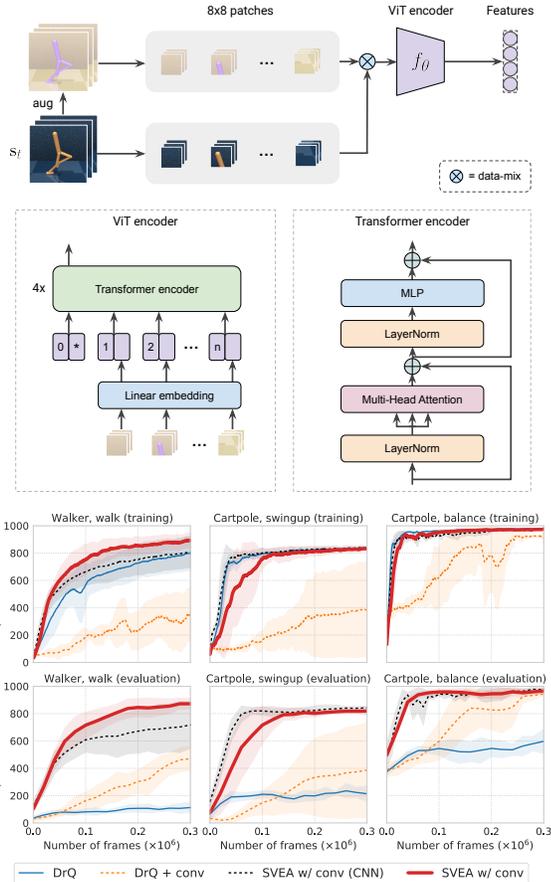


Figure 7. (Top) **ViT architecture.** Observations are divided into 144 non-overlapping space-time patches and linearly projected into tokens. Each token uses a learned positional encoding and we also use a learnable class token as in [10]. The ViT encoder consists of 4 stacked Transformer encoders [69]. (Bottom) **RL with a ViT encoder.** Training and test performance of SVEA and DrQ using ViT encoders. We report results for three tasks and test performance is evaluated on the `color_hard` benchmark of DMControl-GB. Mean of 5 seeds, shaded area is  $\pm 1$  std. deviation.

Table 2. **Generalization in robotic manipulation.** Task success rates of SVEA and DrQ with CNN and ViT encoders in the training environment, as well as aggregated success rates across 25 different test environments with randomized camera pose, colors, lighting, and background. Mean of 5 seeds.

Robotic manipulation	Arch. (encoder)	reach (train)	reach (test)	mv. tgt. (train)	mv. tgt. (test)	push (train)	push (test)
DrQ	CNN	<b>1.00</b>	0.60	<b>1.00</b>	0.69	<b>0.76</b>	0.26
DrQ + conv	CNN	0.59	0.77	0.60	0.89	0.13	0.12
SVEA w/ conv	CNN	<b>1.00</b>	<b>0.89</b>	<b>1.00</b>	<b>0.96</b>	0.72	<b>0.48</b>
DrQ	ViT	0.93	0.14	<b>1.00</b>	0.16	0.73	0.05
DrQ + conv	ViT	0.26	0.67	0.48	<b>0.82</b>	0.08	0.07
SVEA w/ conv	ViT	<b>0.98</b>	<b>0.71</b>	<b>1.00</b>	0.81	<b>0.82</b>	<b>0.17</b>

### 6.3 Robotic Manipulation

We additionally consider a set of goal-conditioned robotic manipulation tasks using a simulated Kinova Gen3 arm: (i) *reach*, a task in which the robot needs to position its gripper above a goal indicated by a red mark; (ii) *reach moving target*, a task similar to (i) but where the robot needs to follow a red mark moving continuously in a zig-zag pattern at a random velocity; and (iii) *push*, a task in which the robot needs to push a cube to a red mark. The initial configuration of gripper, object, and goal is randomized, the agent uses 2D positional control, and policies are trained using dense rewards. Observations are stacks of RGB frames with no access to state information. Training and test environments are shown in Figure 8. See Appendix G for further details and environment samples.

Results are shown in Figure 9 and Figure 10. For both CNN and ViT encoders, SVEA trained with *conv* augmentation has similar sample efficiency and training performance as DrQ trained *without* augmentation, while DrQ + conv exhibits poor sample efficiency and fails to solve the *push* task. Generalization results are shown in Table 2. We find that naïve application of data augmentation has a higher success rate in test environments than DrQ, despite being less successful in the training environment, which we conjecture is because it is optimized only from augmented data. Conversely, SVEA achieves high success rates during both training and testing.

**Conclusion.** SVEA is found to greatly improve both stability and sample efficiency under augmentation, while achieving competitive generalization results. Our experiments indicate that our method scales to ViT-based architectures, and it may therefore be a promising technique for large-scale RL experiments where data augmentation is expected to play an increasingly important role.

**Broader Impact.** While our contribution aims to reduce computational cost of image-based RL, we remain concerned about the growing ecological and economical footprint of deep learning – and RL in particular – with increasingly large models such as ViT; see Appendix J for further discussion.

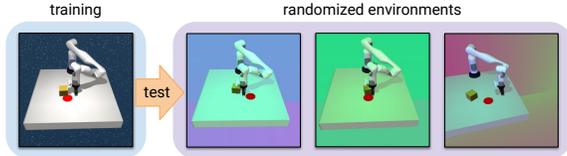


Figure 8. **Robotic manipulation.** Agents are trained in a fixed environment and evaluated on challenging environments with randomized colors, lighting, background, and camera pose.

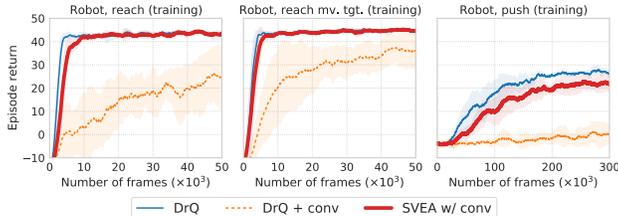


Figure 9. **Stability with a CNN encoder.** Training performance (episode return) of SVEA and DrQ in 3 robotic manipulation tasks. Mean and std. deviation of 5 seeds. Success rates are shown in Table 2.

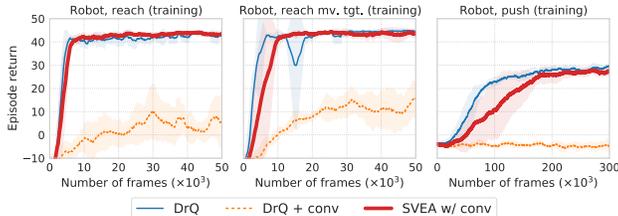


Figure 10. **Stability with a ViT encoder.** Training performance (episode return) of SVEA and DrQ in 3 robotic manipulation tasks. Mean and std. deviation of 5 seeds. Success rates are shown in Table 2. DrQ is especially unstable under augmentation when using a ViT encoder.

**Acknowledgments and Funding Transparency.** This work was supported by grants from DARPA LwLL, NSF CCF-2112665 (TILOS), NSF 1730158 CI-New: Cognitive Hardware and Software Ecosystem Community Infrastructure (CHASE-CI), NSF ACI-1541349 CC\*DNI Pacific Research Platform, NSF grant IIS-1763278, NSF CCF-2112665 (TILOS), as well as gifts from Qualcomm, TuSimple and Picsart.

## References

- [1] Rishabh Agarwal, Marlos C. Machado, P. S. Castro, and Marc G. Bellemare. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. *ArXiv*, abs/2101.05265, 2021.
- [2] Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
- [3] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. Association for Computing Machinery, 2021.
- [4] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Denison, David Farhi, Quirin Fischer, et al. Dota 2 with large scale deep reinforcement learning. *ArXiv*, abs/1912.06680, 2019.
- [5] T. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, J. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, et al. Language models are few-shot learners. *arXiv*, abs/2005.14165, 2020.
- [6] Yevgen Chebotar, A. Handa, Viktor Makoviychuk, M. Macklin, J. Issac, Nathan D. Ratliff, and D. Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. *2019 International Conference on Robotics and Automation (ICRA)*, pages 8973–8979, 2019.
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020.
- [8] K. Cobbe, Oleg Klimov, Christopher Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In *Icml*, 2019.
- [9] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning, 2018.
- [10] A. Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, M. Dehghani, Matthias Minderer, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929, 2020.
- [11] Frederik Ebert, Chelsea Finn, Sudeep Dasari, Annie Xie, Alex Lee, and Sergey Levine. Visual foresight: Model-based deep reinforcement learning for vision-based robotic control, 2018.
- [12] Lasse Espeholt, Hubert Soyer, R. Munos, K. Simonyan, V. Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. *ArXiv*, abs/1802.01561, 2018.
- [13] Jesse Farebrother, Marlos C. Machado, and Michael H. Bowling. Generalization and regularization in dqn. *ArXiv*, abs/1810.00123, 2018.
- [14] Meire Fortunato, M. G. Azar, Bilal Piot, Jacob Menick, Ian Osband, A. Graves, Vlad Mnih, R. Munos, et al. Noisy networks for exploration. *ArXiv*, abs/1706.10295, 2018.
- [15] Scott Fujimoto, H. V. Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. *ArXiv*, abs/1802.09477, 2018.
- [16] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations, 2018.
- [17] R. Givan, T. Dean, and M. Greig. Equivalence notions and model minimization in markov decision processes. *Artificial Intelligence*, 147:163–223, 2003.
- [18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- [19] Danijar Hafner, T. Lillicrap, Ian S. Fischer, R. Villegas, David R Ha, H. Lee, and James Davidson. Learning latent dynamics for planning from pixels. *ArXiv*, abs/1811.04551, 2019.
- [20] Danijar Hafner, T. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *ArXiv*, abs/1912.01603, 2020.
- [21] Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data augmentation. In *International Conference on Robotics and Automation*, 2021.

- [22] Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà, Pieter Abbeel, Alexei A. Efros, Lerrel Pinto, and Xiaolong Wang. Self-supervised policy adaptation during deployment. In *International Conference on Learning Representations*, 2021.
- [23] H. V. Hasselt, A. Guez, Matteo Hessel, V. Mnih, and D. Silver. Learning values across many orders of magnitude. In *Nips*, 2016.
- [24] H. V. Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q-learning. In *Aaai*, 2016.
- [25] M. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In *AAAI Fall Symposia*, 2015.
- [26] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [27] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9726–9735, 2020.
- [28] Matteo Hessel, Joseph Modayil, H. V. Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, et al. Rainbow: Combining improvements in deep reinforcement learning. In *Aaai*, 2018.
- [29] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks, 2016.
- [30] Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 1998.
- [31] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *ArXiv*, abs/1806.10293, 2018.
- [32] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.
- [33] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *International Conference on Learning Representations*, 2020.
- [34] K. G. Larsen and A. Skou. Bisimulation through probabilistic testing (preliminary report). In *Proceedings of the 16th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1989.
- [35] Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *arXiv preprint arXiv:2004.14990*, 2020.
- [36] Kimin Lee, Kibok Lee, Jinwoo Shin, and Honglak Lee. A simple randomization technique for generalization in deep reinforcement learning. *ArXiv*, abs/1910.05396, 2019.
- [37] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [38] T. Lillicrap, J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *CoRR*, abs/1509.02971, 2016.
- [39] L. J. Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8:293–321, 2004.
- [40] Xingyu Lin, Harjatin Singh Baweja, George Kantor, and David Held. Adaptive auxiliary task weighting for reinforcement learning. *Advances in neural information processing systems*, 32, 2019.
- [41] Clare Lyle, Mark Rowland, Georg Ostrovski, and Will Dabney. On the effect of auxiliary tasks on representation dynamics. In *Aistats*, 2021.
- [42] P. Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andy Ballard, Andrea Banino, Misha Denil, R. Goroshin, et al. Learning to navigate in complex environments. *ArXiv*, abs/1611.03673, 2017.
- [43] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [44] Ashvin V Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforcement learning with imagined goals. In *Advances in Neural Information Processing Systems*, pages 9191–9200, 2018.
- [45] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016.
- [46] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016.

- [47] Deepak Pathak, Dhiraj Gandhi, and A. Gupta. Self-supervised exploration via disagreement. *ArXiv*, abs/1906.04161, 2019.
- [48] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. *2018 IEEE International Conference on Robotics and Automation (ICRA)*, May 2018.
- [49] Lerrel Pinto and Abhinav Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *2016 IEEE international conference on robotics and automation (ICRA)*, pages 3406–3413. Ieee, 2016.
- [50] Lerrel Pinto, Marcin Andrychowicz, Peter Welinder, Wojciech Zaremba, and Pieter Abbeel. Asymmetric actor critic for image-based robot learning. *arXiv preprint arXiv:1710.06542*, 2017.
- [51] Roberta Raileanu, M. Goldstein, Denis Yarats, Ilya Kostrikov, and R. Fergus. Automatic data augmentation for generalization in deep reinforcement learning. *ArXiv*, abs/2006.12862, 2020.
- [52] Fabio Ramos, Rafael Possas, and Dieter Fox. Bayessim: Adaptive domain randomization via probabilistic inference for robotics simulators. *Robotics: Science and Systems XV*, Jun 2019.
- [53] Tom Schaul, John Quan, Ioannis Antonoglou, and D. Silver. Prioritized experience replay. *CoRR*, abs/1511.05952, 2016.
- [54] Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bachman. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint arXiv:2007.05929*, 2020.
- [55] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak. Planning to explore via self-supervised world models, 2020.
- [56] Evan Shelhamer, Parsa Mahmoudieh, Max Argus, and Trevor Darrell. Loss is its own reward: Self-supervision for reinforcement learning. *ArXiv*, abs/1612.07307, 2017.
- [57] Connor Shorten and T. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:1–48, 2019.
- [58] Xingyou Song, Yiding Jiang, Stephen Tu, Yilun Du, and Behnam Neyshabur. Observational overfitting in reinforcement learning. *ArXiv*, abs/1912.02975, 2020.
- [59] Aravind Srinivas, Michael Laskin, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv preprint arXiv:2004.04136*, 2020.
- [60] Austin Stone, Oscar Ramirez, K. Konolige, and Rico Jonschkowski. The distracting control suite - a challenging benchmark for reinforcement learning from pixels. *ArXiv*, abs/2101.02722, 2021.
- [61] Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling representation learning from reinforcement learning. *ArXiv*, abs/2004.1499, 2020.
- [62] R. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9–44, 2005.
- [63] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- [64] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, et al. Deepmind control suite. Technical report, DeepMind, 2018.
- [65] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- [66] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep 2017.
- [67] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012. doi: 10.1109/iros.2012.6386109.
- [68] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2018.
- [69] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv*, abs/1706.03762, 2017.
- [70] Oriol Vinyals, I. Babuschkin, Wojciech Czarnecki, Micha "el Mathieu, Andrew Dudzik, J. Chung, D. Choi, Richard Powell, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, pages 1–5, 2019.
- [71] K. Wang, Bingyi Kang, Jie Shao, and Jiashi Feng. Improving generalization in reinforcement learning with mixture regularization. *ArXiv*, abs/2010.10814, 2020.

- [72] Xudong Wang, Long Lian, and Stella X. Yu. Unsupervised visual attention and invariance for reinforcement learning. *ArXiv*, abs/2104.02921, 2021.
- [73] Ziyu Wang, Tom Schaul, Matteo Hessel, H. V. Hasselt, Marc Lanctot, and N. D. Freitas. Dueling network architectures for deep reinforcement learning. *ArXiv*, abs/1511.06581, 2016.
- [74] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.
- [75] Zhenlin Xu, Deyi Liu, Junlin Yang, and M. Niethammer. Robust and generalizable visual representation learning via random convolutions. *ArXiv*, abs/2007.13003, 2020.
- [76] Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images, 2019.
- [77] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with prototypical representations. *arXiv preprint arXiv:2102.11271*, 2021.
- [78] Sarah Young, Dhiraj Gandhi, Shubham Tulsiani, Abhinav Gupta, Pieter Abbeel, and Lerrel Pinto. Visual imitation made easy. *CoRL*, 2020.
- [79] Tianhe Yu, Saurabh Kumar, A. Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *ArXiv*, abs/2001.06782, 2020.
- [80] A. Zhang, Rowan McAllister, R. Calandra, Y. Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. *ArXiv*, abs/2006.10742, 2020.
- [81] C. Zhang, Oriol Vinyals, R. Munos, and S. Bengio. A study on overfitting in deep reinforcement learning. *ArXiv*, abs/1804.06893, 2018.
- [82] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016.
- [83] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *Icra*, pages 3357–3364. Ieee, 2017.
- [84] Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *Proceedings of the 23rd National Conference on Artificial Intelligence*, volume 3, 2008.