

# Approximate Equivariance in Reinforcement Learning

Jung Yeon Park\*

Alec Koppel†

\*Northeastern University

Sujay Bhatt†

Sumitra Ganesh†

Sihan Zeng†

†J.P.Morgan AI Research

Lawson L.S. Wong\*

Robin Walters\*

## Abstract

Equivariant neural networks have shown great success in reinforcement learning, improving sample efficiency and generalization when there is symmetry in the task. However, in many problems, only approximate symmetry is present, which makes imposing exact symmetry inappropriate. Recently, approximately equivariant networks have been proposed for supervised classification and modeling physical systems. In this work, we develop approximately equivariant algorithms in reinforcement learning (RL). We define approximately equivariant MDPs and theoretically characterize the effect of approximate equivariance on the optimal  $Q$  function. We propose novel RL architectures using relaxed group and steerable convolutions and experiment on several continuous control domains and stock trading with real financial data. Our results demonstrate that the approximately equivariant network performs on par with exactly equivariant networks when exact symmetries are present, and outperforms them when the domains exhibit approximate symmetry. As an added byproduct of these techniques, we observe increased robustness to noise at test time. Our code is available at [https://github.com/jypark0/approx\\_equiv\\_rl](https://github.com/jypark0/approx_equiv_rl).

## 1 INTRODUCTION

Symmetry is a powerful inductive bias that can be used to improve generalization and data efficiency in deep learning. One way to leverage symmetry

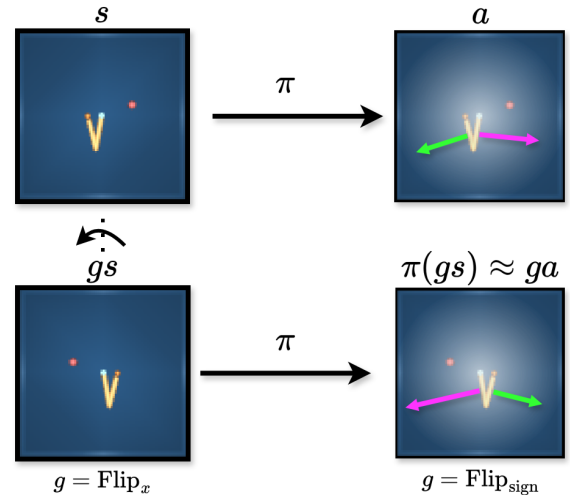


Figure 1: An approximately equivariant policy  $\pi$  on a Reacher domain, where the goal is to determine the torques (green, magenta) to apply on each joint for the fingertip to reach the target (red). Due to wear, the first joint is more responsive to positive torques. When the state is flipped, the policy also flips the actions but can learn to adjust for symmetry breaking factors.

is through equivariant neural networks, which are model classes constrained to respect the symmetry of a known ground truth. Equivariant neural networks have successfully been applied to image classification (Cohen and Welling, 2016; Worrall et al., 2017), particle physics (Bogatskiy et al., 2020), molecular biology (Satorras et al., 2021; Thomas et al., 2018), and robotic manipulation (Wang et al., 2022b). Empirical studies have demonstrated that equivariant networks require much fewer data than their standard network counterparts (Winkels and Cohen, 2018; Wang et al., 2022b), can have fewer parameters (Weiler and Cesa, 2019; He et al., 2022), and can generalize better to unseen data (Wang et al., 2020; Fuchs et al., 2020).

However, equivariant neural networks crucially assume that the data is perfectly symmetric in both the inputs and outputs, which may not be true in real-world data such as fluid dynamics (Wang et al., 2022c) or finan-

cial data (Black, 1986). By relaxing the strict equivariance constraints, approximately equivariant networks can outperform exactly equivariant and unconstrained networks in the presence of asymmetry. While various approaches to achieve approximate equivariance have been proposed (Wang et al., 2022c; van der Ouderaa et al., 2022; McNeela, 2023; Kim et al., 2023), they focus on vision-based tasks or dynamics modeling.

One area where symmetry has been especially useful is in reinforcement learning (RL), where equivariant networks greatly improve sample efficiency (Wang et al., 2022b; Zhu et al., 2022), a key challenge in RL. However, most works consider exact symmetry and use exactly equivariant networks, which cannot address symmetry breaking in the reward or transition functions or noise in the observations. In this work, we employ relaxed group and steerable convolutional neural networks for RL (Wang et al., 2022c); they are flexible enough to adapt to approximate equivariance but also have improved efficiency and robustness.

In this paper, we theoretically and empirically investigate approximately equivariant reinforcement learning. Our key contributions are to:

- formalize the notion of approximately equivariant MDPs and prove the (optimal) value function in such MDPs exhibits approximate equivariance, motivating the use of approximately equivariant RL,
- introduce a novel approximately equivariant RL architecture using relaxed group convolutions,
- demonstrate improved sample efficiency and robustness to noise for our approximately equivariant RL compared to other baselines with or without symmetry biases,
- successfully apply approximate equivariant RL to real-world financial data.

## 2 RELATED WORK

**Equivariant Reinforcement Learning** Early works explored equivalence classes in reinforcement learning from the lens of abstractions by defining MDP homomorphisms (Ravindran and Barto, 2002; Zinkevich and Balch, 2001). More recently, several approaches have combined function approximation with RL with equivariant neural networks (Van der Pol et al., 2020; Wang et al., 2022b; Mondal et al., 2020) with significantly improved sample efficiency. However, all of these works considered perfectly symmetric domains where the policy is constrained to be exactly equivariant. This paper considers domains with symmetry breaking factors where exactly equivariant networks can be suboptimal.

**Approximate Equivariant Architectures** There has been recent interest in exploring approximate equivariance and approximately equivariant neural networks (Finzi et al., 2021; Wang et al., 2022c; Romero and Lohit, 2022; van der Ouderaa et al., 2022; McNeela, 2023; Petrache and Trivedi, 2024; Samudre et al., 2024). Wang et al. (2022c, 2024b) use a linear combination of exactly equivariant convolution kernels with learnable weights to achieve relaxed equivariance and discover symmetry breaking factors. van der Ouderaa et al. (2022) define a nonstationary kernel and a tunable frequency parameter to control the amount of approximate equivariance. McNeela (2023) propose using a neural network to approximate the exponential map from the Lie algebra to the group to learn almost equivariant functions. Petrache and Trivedi (2024) give theoretical bounds on when approximate equivariance can improve generalization. However, none of these works studied approximate equivariance in RL, the main focus of this work.

Closest to our setting is Residual Pathway Priors (Finzi et al., 2021), which considered soft equivariance constraints in model-free RL. They construct a relaxed equivariant neural network layer as the sum of an exactly equivariant and a non-equivariant layer with a prior on the equivariant layer. We take a different approach in this work and use relaxed group convolutions Wang et al. (2022c), which are flexible enough to learn different outputs for each transformation.

**Learning with Latent Symmetry** Other works also apply equivariant neural networks to domains with latent symmetry. These are cases where the full state has exact symmetry but only partial observations with an unknown group action are available to the model. Park et al. (2022) learn the out-of-plane rotations from 2D images using a symmetric embedding network while others have learned 3D rotational features from images using manifold latent variables (Falorsi et al., 2018) or disentanglement (Quessard et al., 2020). Wang et al. (2022a) find that equivariant models where the group acts directly on observation space perform well in RL even with camera skew or occlusions. They define extrinsic equivariance (transformed samples are outside the data distribution) and show that it can benefit in some scenarios but can also be harmful in certain cases (Wang et al., 2024a). Unlike these works where the observation is partial and does not contain full information about the state, we assume that the domains are fully observable and consider various symmetry breaking factors.

### 3 BACKGROUND

In this section, we provide some background on symmetry groups and equivariant functions. As building blocks of exactly and approximately equivariant networks, we also describe exact and relaxed group convolutions, respectively.

#### 3.1 Groups and Equivariance

A symmetry group  $G$  is a set equipped with a binary operation that satisfies associativity, existence of an identity, and existence of inverses. A group can act on vector space  $X$  via a group representation  $\rho_X$  which homomorphically assigns each element  $g \in G$  an invertible matrix  $\rho_X(g) \in \text{GL}(X)$ . For example, for a finite group  $G$ , the regular representation acts on  $\mathbb{R}^{|G|}$  by permuting basis elements  $\{e_g : g \in G\}$  as  $\rho_{\text{reg}}(h)e_g = e_{hg}$ . A function  $f : X \rightarrow Y$ ,  $x \mapsto y$  is  $G$ -equivariant if  $f(\rho_X(g)(x)) = \rho_Y(g)f(x)$ . That is, transformations of the input  $x$  by  $g$  correspond to transformations of the output by the same group element. We can enforce this constraint in equivariant neural networks to learn only over the space of equivariant functions by replacing linear layers with group or steerable convolutional layers. One benefit of enforcing equivariance is lower sample complexity as the network searches over a reduced function class.

#### 3.2 Group Convolution

One method of constructing equivariant network layers is by group convolution (Cohen and Welling, 2016), which we briefly describe here. Group convolutions map between features which are signals over the group  $f : G \rightarrow \mathbb{R}$ . For inputs not natively of this form, a lift operation must first be performed. Let  $\psi_\theta : G \rightarrow \mathbb{R}$  be the convolutional kernel parameterized by  $\theta$ . A  $G$ -equivariant group convolutional layer is defined as

$$(f \star \psi_\theta)(g) = \sum_{h \in G} \psi_\theta(g^{-1}h)f(h). \quad (1)$$

Equivariance follows from the fact that the kernel depends only on the product  $g^{-1}h$  and not the specific elements  $(g, h)$ . For example, if we consider equivariance across translations, we obtain the standard convolution where  $h, g \in \mathbb{Z}^2$  and  $g^{-1}h = h - g$ . Another possible approach to constructing equivariant network layers is with  $G$ -steerable convolutions (Cohen and Welling, 2017), which can generalize to continuous groups.

#### 3.3 Relaxed Group Convolution

A key component of our method is the relaxed version of the group convolution (Wang et al., 2022c). The

kernel  $\psi$  is replaced with several kernels  $\{\psi_l\}_{l=1}^L$  and the output is composed as a linear combination. The relaxed group convolution is defined as

$$(f \widetilde{\star} \psi_\theta)(g) = \sum_{h \in G} f(h) \sum_{l=1}^L w^l(h) \psi_\theta^l(g^{-1}h), \quad (2)$$

where  $w^l$  are the relaxed weights and each  $\psi_\theta^l$  are constrained to be exactly equivariant. Note that as  $w^l(h)$  depends on the specific element  $h$ , this breaks the strict equivariance of the group convolution. Wang et al. (2022c) also introduce relaxed versions of steerable convolutions, see Wang et al. (2022c) or Appendix C.2 for more details.

#### 3.4 Approximate Equivariance

There have been several different definitions of approximate, relaxed, or partial equivariance. In this paper, we use the definition given by Petrache and Trivedi (2024). We give some background to build up to the definition. Let  $G$  be a group and  $f : X \rightarrow Y$ ,  $x \mapsto y$  be the task function.

**Definition 1** (Equivariance Error). *For  $g \in G$  and  $x \in X$ , the equivariance error  $ee(f, g, x)$  is defined as*

$$ee(f, g, x) = \|f(g(x)) - g(f(x))\|, \quad (3)$$

Equivariance error measures exactly how far a function is from perfect equivariance with respect to  $G$  for a particular  $x$ . For an exactly  $G$ -equivariant function,  $ee(f, g, x) = 0$  for all  $g \in G$  and  $x \in X$ .

**Definition 2** ( $\varepsilon$ -stabilizer). *The  $\varepsilon$ -stabilizer of  $f$  and  $G$  is defined as*

$$\text{Stab}_\varepsilon(f, G) = \{g \in G \mid ee(f, g, x) \leq \varepsilon\}. \quad (4)$$

The  $\varepsilon$ -stabilizer gives the set of group elements for which the equivariance error is under some threshold.

**Definition 3** (Approximate  $G$ -Equivariance). *Given a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  and a group  $G$ ,  $f$  is approximately  $G$ -equivariant if  $\text{Stab}_\varepsilon(f, G) = G$ .*

We adopt the definition of approximate equivariance where  $f$  has bounded equivariance error for all  $g \in G$ , in contrast to *partial equivariance*, where  $\text{Stab}_\varepsilon(f, G) < G$ .

## 4 METHOD: APPROXIMATELY EQUIVARIANT REINFORCEMENT LEARNING

We first theoretically characterize the problem by defining approximately equivariant Markov decision

processes (MDP). We then prove that environments with approximate symmetry admit approximately invariant  $Q$  functions. This motivates our method of using approximately equivariant neural networks to learn the policy and  $Q$  function.

#### 4.1 Approximately Equivariant MDP

Consider an infinite-horizon discounted-reward Markov decision process (MDP) represented by a tuple  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$  with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , instantaneous reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , a transition function  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  and discount factor  $\gamma \in (0, 1)$ .

Let  $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  be a policy giving the probability  $\pi(a|s)$  of taking action  $a$  in state  $s$ . The expected cumulative reward of using the policy starting from state  $s$  (or state  $s$  and action  $a$ ) are the value functions defined as follows

$$\begin{aligned} V^\pi(s) &:= \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, a_k) \middle| s_0 = s \right], \\ Q^\pi(s, a) &:= \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, a_k) \middle| s_0 = s, a_0 = a \right]. \end{aligned} \quad (5)$$

The goal is to find a policy  $\pi^*$  that maximizes the expected return with an initial state distribution  $\xi$

$$\pi^* := \arg \max_{\pi} \mathbb{E}_{s_0 \sim \xi} [V^\pi(s_0)].$$

We denote  $V^* = V^{\pi^*}$  and  $Q^* = Q^{\pi^*}$ .

Let  $G$  be a group acting on  $\mathcal{S}$  and  $\mathcal{A}$ . Denote the action of an element  $g \in G$  on  $s$  and  $a$  by  $gs$  and  $ga$ , respectively. We now extend the definition of Equivariant MDPs (Van der Pol et al., 2020) to cases where the symmetry is approximate.

**Definition 4.** An MDP is  $(G, \epsilon_R, \epsilon_P)$ -invariant if

$$\begin{aligned} |R(gs, ga) - R(s, a)| &\leq \epsilon_R, \forall g \in G \\ d_{\mathcal{F}}(P(gs' | gs, ga), P(s' | s, a)) &\leq \epsilon_P, \forall g \in G, \end{aligned}$$

where  $d_{\mathcal{F}}(\mu, \nu) := \sup_{f \in \mathcal{F}} \left| \int_{\mathcal{S}} f d\mu - \int_{\mathcal{S}} f d\nu \right|$  is an integral probability metric (IPM) between two distributions  $\mu, \nu \in \Delta(\mathcal{X})$ .

Some well known examples of IPM include (Sriperumbudur et al., 2009): total variation distance ( $\mathcal{F} = \{f : \|f\|_{\infty} \leq 1\}$ ) and Kantorovich metric ( $\mathcal{F} = \{f : \|f\|_{\text{Lip}} \leq 1\}$ ). A useful property of IPMs is, given a function class  $\mathcal{F}$  and a function  $f$  (Müller, 1997)

$$\left| \int_{\mathcal{S}} f d\mu - \int_{\mathcal{S}} f d\nu \right| \leq \rho_{\mathcal{F}}(f) \cdot d_{\mathcal{F}}(\mu, \nu),$$

where the Minkowski functional w.r.t  $\mathcal{F}$  is

$$\rho_{\mathcal{F}}(f) = \inf \{ \rho \in \mathbb{R}_{\geq 0} : \rho^{-1} f \in \mathcal{F} \}.$$

For the total variation distance  $\rho_{\mathcal{F}}(f) := \frac{1}{2}(\max f - \min f)$  and for Kantorovich metric  $\rho_{\mathcal{F}}(f) := \|f\|_{\text{Lip}}$ .

The following theorem provides a characterization of the gap between the value functions in the original and symmetry transformed domain, for the  $(G, \epsilon_R, \epsilon_P)$ -invariant MDP described in Definition 4. Theorem 1 highlights that the  $Q$ -function is approximately group-invariant, where the approximation is now a function of the reward and transition mismatch, the discount factor, and the Minkowski functional evaluated on the optimal value function.

**Theorem 1.** Let the rewards  $R$  be bounded  $R_{\min} \leq R \leq R_{\max}$ ,  $0 \leq \gamma < 1$  and let  $g \in G$  be an onto mapping. For any state  $s$  and action  $a$ , we have

$$\begin{aligned} |Q^*(s, a) - Q^*(gs, ga)| &\leq \alpha, \\ |V^*(s) - V^*(gs)| &\leq \alpha, \end{aligned}$$

where  $\alpha = \frac{\epsilon_R + \gamma \rho_{\mathcal{F}}(V^*) \epsilon_P}{1 - \gamma}$ .

Theorem 1 implies that when the invariance mismatch is small – i.e., when the domain has only minor symmetry violations – the  $Q$ -function is approximately group-invariant. A proof is provided in Appendix A. Note that, in Theorem 1, when the Kantorovich metric is used for uncertainty characterization,  $\rho_{\mathcal{F}}(V^*) = \|V^*\|_{\text{Lip}}$ , where  $\|\cdot\|_{\text{Lip}}$  is the Lipschitz norm of the value function (Gelada et al., 2019). For total variation distance,  $\rho_{\mathcal{F}}(V^*) = |R_{\max} - R_{\min}|$ .

Also, from Theorem 1, it is clear that when  $\gamma \in [0, 1)$ , we obtain a non-trivial characterization, while  $\gamma = 1$  results in a trivial and uninformative bound. This is the limitation of the infinite horizon setting, and can be remedied by considering an arbitrary finite-horizon setup. We do this for the sake of completeness in Appendix B. We not only show that the finite horizon setup allows for a time-dependent transition function, but also obtain an approximate group-invariance of the time dependent  $Q$ -function in terms of similar elements that appear in Theorem 1.

There are different ways to use the above results. One can discover how approximate the value functions are and learn  $\alpha$ , or one can incorporate approximate equivariance into the model and leverage the benefits of equivariance. We take the latter approach and consider approximately equivariant networks for the policy and critic in domains with inexact symmetry.

#### 4.2 Approximately Equivariant Actor-Critic

We propose approximately equivariant versions of two commonly used actor-critic algorithms, DrQv2 (Yarats



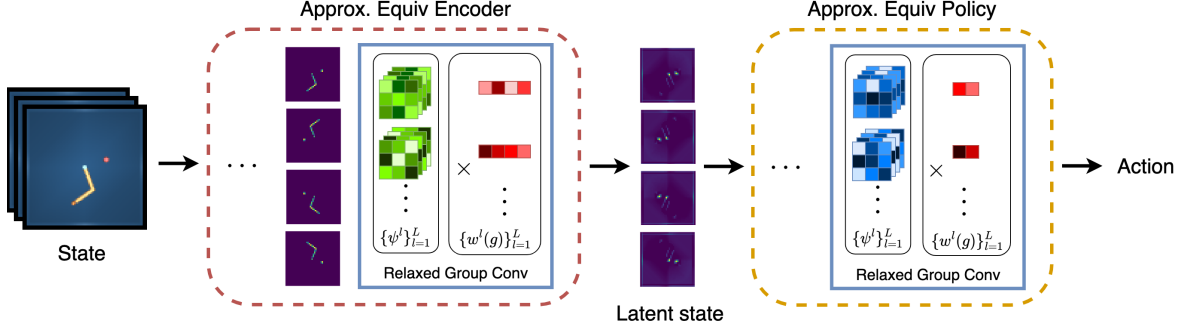


Figure 2: Illustration of the approximately  $D_2$ -equivariant encoder and policy (critic is not shown for space). The  $D_2$  group consists of vertical reflections and  $\pi$  rotations. Both the encoder and policy consist of relaxed group convolution layers.

et al., 2021) and SAC (Haarnoja et al., 2018). In doing so, we generalize exactly equivariant versions of SAC (Wang et al., 2022b) and DrQv2 (Wang et al., 2022a) from previous works by replacing strictly equivariant layers with relaxed equivariant layers.

**Illustrative Example** We first illustrate how to apply our proposed approximately equivariant actor-critic architecture on the **Reacher** domain; see Figure 2. The objective is to actuate a two-joint arm so that the end effector reaches the red point. The state is a stack of consecutive images  $s \in \mathbb{R}^{C \times H \times W}$  and the action  $a \in \mathbb{R}^2$  corresponds to torques for the first and second arms. There is clear rotational and reflectional symmetry in this domain. If the state (image) is rotated, the action should be invariant to rotations as they are angular torques. If the state is reflected, then the action would also correspondingly be flipped (in sign). However, as in the example in Figure 1, the first joint is more responsive to positive torques, which breaks rotational and reflectional symmetry.

For this domain, we implement approximate equivariance to the group  $D_2$  of vertical reflections and  $\pi$  rotations. The group  $D_2$  transforms the input states by image transformations, where the input images are reflected or rotated. Latent representations are images  $z: \mathbb{R}^2 \rightarrow \mathbb{R}^C$  where  $g \in D_2$  acts on the pixel axes by image transformation and on the channel axis by permutations corresponding to the regular representation of  $D_2$ , i.e.  $(gz)(x, y) = \rho_{\text{reg}}(g)z(g^{-1} \cdot (x, y))$ . Note that the latent representations can be high-dimensional, consisting of a direct sum of several different or repeated low-dimensional representations of  $D_2$ . For the output, the torques  $a_1$  and  $a_2$  are scalars that change sign under reflection but are invariant under rotations.

**Encoder, Policy, and Critic** We extend exactly equivariant versions of SAC (Wang et al., 2022b) and DrQv2 (Wang et al., 2022a) by replacing each group convolution with relaxed group convolutions for the

encoder, policy, and critics. Practically, each relaxed group convolution layer contains  $L$  exactly equivariant kernels  $\psi_l$  and the output is a linear combination of the outputs of these convolutions and relaxed weights  $w^l(g)$ . The  $w^l(g)$  also transform as the regular representation of  $G$ , see Section 3.1 for the definition for finite groups.

The encoder  $E$  and the policy  $\pi$  are approximately equivariant. The latent state  $z$  output by  $E$  is defined to transform as the direct sum of regular representations of  $G$ . The action representation is domain-specific. The critics are approximately invariant and output scalars  $q_{(s,a)}$  that are fixed by  $G$ , i.e. transform via the trivial representation. For more details, please see Section 5 and Appendix C.

In the case of continuous groups, we can also construct relaxed steerable versions of the encoder, policy, and critics. Analogous to the group convolution case, we can replace the exactly equivariant steerable convolutions with relaxed steerable convolutions. See Appendix C for more details.

## 5 EXPERIMENTS

We experiment on how approximately equivariant RL compares to methods with exact equivariance and no equivariance in domains with both exact symmetry and various symmetry breaking factors, and to elucidate when approximate equivariance should be preferred. We consider standard continuous control domains and stock trading with real-world data.

### 5.1 Continuous Control

We first experiment on four continuous control domains in DeepMind Control Suite (Tassa et al., 2018). Similar to Wang et al. (2022a), we consider a subset of the domains which have apparent symmetry. **Acrobot**, **Cartpole**, and **BallInCup** have reflectional symmetry

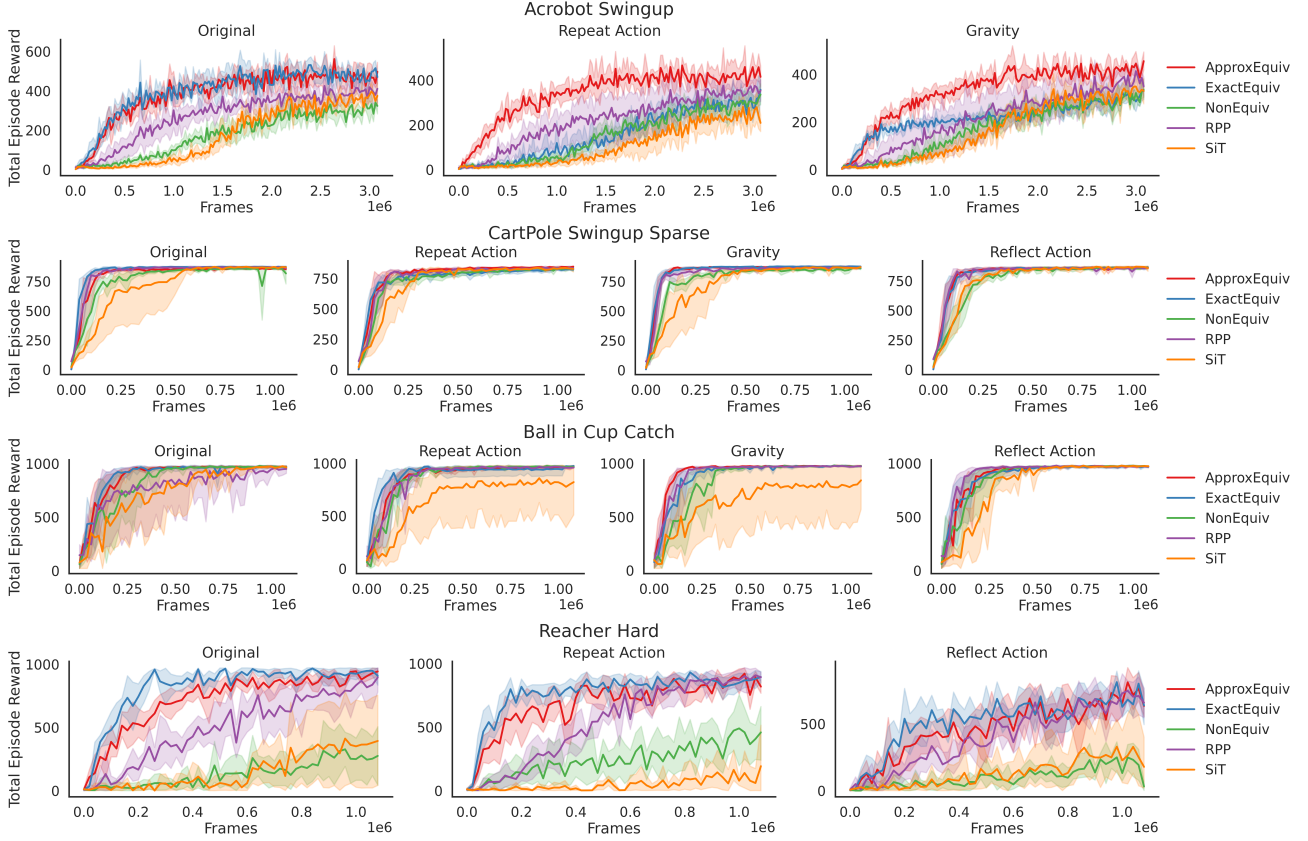


Figure 3: Total episode reward on selected domains in the DeepMind Control Suite, shaded regions indicate 95% confidence intervals (CI). Compared to an exactly equivariant agent (ExactEquiv), our approximately equivariant agent (ApproxEquiv) outperforms in **Acrobot**, performs similarly in two domains, and is slightly worse in the **Reacher** domain. ApproxEquiv can outperform ExactEquiv on some modified variants with inexact symmetry as it can adjust for symmetry breaking. Our agent outperforms all other baselines, including a non-equivariant agent, suggesting that relaxed symmetry is a good inductive bias.



Figure 4: Selected domains in DeepMind Control Suite. The domains were modified to remove extrinsic symmetry and to include several types of symmetry breaking factors such as repeating or reflecting actions in certain states, or by modifying gravity.

described by the group  $D_1$  and **Reacher** has  $D_2$  symmetry. For all domains, the observations are a stack of 3 consecutive RGB images.

We modify the domains to carefully control the type and degree of symmetry breaking that is present. We first remove fixed background features such as random stars in the sky and checkered floors (see Figure 4). These features break symmetry to some extent since

they do not transform with the underlying state, but give a form of mild symmetry breaking termed extrinsic equivariance, which has an inconsistent impact on equivariant models (Wang et al., 2022a). We then introduce several different symmetry breaking factors for each domain: 1) **repeat\_action** - the action is repeated twice in a certain region of the domain, 2) **gravity** - gravity is modified from the force vector  $(0, 0, -9.81)$  to  $(a, -a, -9.81)$  where  $a \neq 0$ , and 3) **reflect\_action** - the action direction is flipped in certain regions of the domain. **repeat\_action** and **reflect\_action** test local symmetry breaking factors, while **gravity** tests a global symmetry breaking factor. See Appendix D.1 for more details.

**Models** For the continuous control tasks, we implement an approximately equivariant (ApproxEquiv) version of a SOTA image-based RL algorithm DrQv2 (Yarats et al., 2021). We compare with exactly equivariant (ExactEquiv) and non equivariant (NonEquiv)

Table 1: Total episode reward on 50 rollouts for the best policy in the original and noisy domains. Gray values indicate 95% CI. **ApproxEquiv** learns a better policy than baselines on the modified domains and is more robust to noisy inputs.

		No Noise			Noisy		
		ApproxEquiv	ExactEquiv	NonEquiv	ApproxEquiv	ExactEquiv	NonEquiv
ACROBOT	Original	389 $\pm$ 11	<b>522<math>\pm</math>21</b>	309 $\pm$ 22	344 $\pm$ 14	<b>402<math>\pm</math>22</b>	190 $\pm$ 14
	Gravity	<b>471<math>\pm</math>17</b>	382 $\pm$ 15	358 $\pm$ 23	<b>369<math>\pm</math>15</b>	218 $\pm$ 10	202 $\pm$ 12
CARTPOLE	Original	876 $\pm$ 0.2	<b>881<math>\pm</math>0.1</b>	<b>881<math>\pm</math>0.1</b>	778 $\pm$ 22	<b>855<math>\pm</math>0.6</b>	572.5 $\pm$ 25
	Repeat Action	<b>859<math>\pm</math>0.4</b>	749 $\pm$ 13	855 $\pm$ 0.6	<b>624<math>\pm</math>6.0</b>	523 $\pm$ 21	192 $\pm$ 5.0
BALL IN CUP	Original	961 $\pm$ 0.0	958 $\pm$ 0.0	<b>970<math>\pm</math>0.0</b>	<b>882<math>\pm</math>7.7</b>	783 $\pm$ 24	0 $\pm$ 0.0
	Gravity	<b>969<math>\pm</math>0.0</b>	966 $\pm$ 0.0	959 $\pm$ 0.0	<b>888<math>\pm</math>13.2</b>	0 $\pm$ 0.0	1.8 $\pm$ 1.8
REACHER	Original	903 $\pm$ 33	<b>950<math>\pm</math>15</b>	519 $\pm$ 68	<b>778<math>\pm</math>41</b>	745 $\pm$ 44	247 $\pm$ 52
	Reflect Action	<b>757<math>\pm</math>55</b>	707 $\pm$ 59	243 $\pm$ 58	<b>659<math>\pm</math>42</b>	217 $\pm$ 41	82 $\pm$ 29

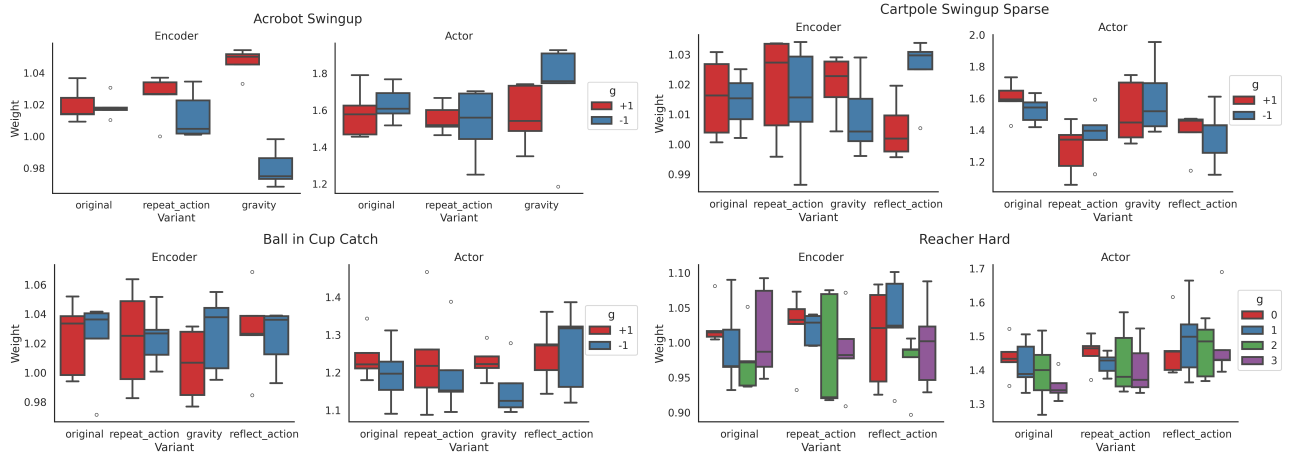


Figure 5: Visualization of relaxed weights for the first layer of the encoder and policy over all runs. Similar weights for each  $g$  indicate perfect equivariance while differing values indicate symmetry breaking. The modified variants of most domains exhibit larger differences or increased variance in the relaxed weights compared to the **original** variant.

versions of the same architecture. We largely use the hyperparameters from [Yarats et al. \(2021\)](#) but reduce the latent dimension for more tractable computation for all methods. We also compare against an approximately equivariant model, Residual Pathway Priors (RPP) ([Finzi et al., 2021](#)), and a self-supervised symmetry-aware model, SiT ([Weissenbacher et al., 2024](#)). We extend RPP to the DrQv2 architecture by using RPP layers in the encoder, policy, and critics. We find that RPP is somewhat sensitive to the speed  $\tau$  of the critic moving average (as mentioned in the original paper), and had to reduce its value for **Acrobot** and **BallInCup** for stability. We also extend SiT to the DrQv2 architecture by using a SiT as the encoder and standard MLPs for the policy and critics. Although we adapted the code from the official SiT implementation, we were unable to modify the input image sizes and had to use the image size used in the original paper (64px).

**Results** Figure 3 show the total episode reward over training. As expected, we confirm that **NonEquiv** has much lower sample efficiency than the models with a symmetry bias. In the **repeat\_action** and **reflect\_action** variants of **Acrobot**, **ApproxEquiv** significantly outperforms **ExactEquiv** and RPP. It does slightly worse than **ExactEquiv** on the **Reacher** domain but beats RPP, suggesting that the symmetry breaking we introduced was not strong enough to achieve incorrect equivariance. It is also possible that **ExactEquiv** can infer the symmetry breaking factors from the 3 frames of input, making the task a case of extrinsic equivariance where an equivariant model can succeed [Wang et al. \(2022a\)](#). In **CartPole** and **BallInCup**, all methods perform similarly and learn an optimal policy quickly. In domains with exact symmetry (**original**), our method **ApproxEquiv** performs similarly to **ExactEquiv**, showing there is no cost in performance by giving the model the ability to adapt

to symmetry breaking in cases where it is not needed. This result supports Proposition 3.1 from Wang et al. (2024b), which proves that relaxed group convolutions initialized to be exactly equivariant stay exactly equivariant when trained with exact data symmetry.

We visualize the relaxed weights of the first layers of the encoder and policy over all runs in Figure 5. If these weights are equal, the model is equivariant; the more they differ the more the model has relaxed the symmetry constraint. For **Acrobot** and **CartPole**, the weights differ more for the modified domains than the original symmetric domain, especially for the encoder, while the policy weights vary more for the modified domains of **BallInCup**. This indicates the relaxed equivariant models have adapted to the symmetry breaking in the domains.

To quantitatively evaluate the models, we select the best-performing policy from all runs and measure the total reward over 50 episodes. The results echo the training curves in Figure 3, where **ApproxEquiv** performs well, particularly in the domains with symmetry breaking factors (see Table 1).

To test whether approximately equivariant models are robust to noisy observations, we also consider variants of the domains where Gaussian noise are added to the input images only at test time ( $\sigma = 0.02$  for **Acrobot** and **Reacher**,  $\sigma = 0.06$  for **CartPole** and **BallInCup**). Interestingly, we find that our approach is more robust to noisy inputs than **ExactEquiv** or **NonEquiv**, especially on the **BallInCup** and **Reacher** domains. We further experiment with *training* on noisy data and test on noisy domains to see which policies are more robust, see Table 4 in Appendix E. We find that in the **BallInCup** domains, the approximately equivariant agent is still more robust to noise than the fully equivariant or non equivariant baselines.

## 5.2 Stock Trading

We also consider a stock trading task using real world price data, formulated as an MDP (Liu et al., 2018). Given a fixed amount of initial cash, the objective is to learn the optimal number of stocks to buy and sell (once daily) to maximize the portfolio value. The state consists of the current cash balance, the stock prices, the number of shares in the current portfolio, and other technical indicators of each stock. The actions are the number of stocks to buy and sell for each stock. The reward is the scaled difference in portfolio values between consecutive timesteps. We assume that the market dynamics are not affected by our trading. There is a small 0.1% transaction cost for every trade. We use real financial data scraped from Yahoo Finance (yfi, 1997) and consider the stocks in the

Table 2: Test results on the stock trading dataset. Gray values indicate 95% CI over 5 runs. The approximately equivariant agents for both scale-translation (ST) and translation (T) outperform the exactly equivariant and non equivariant methods.

		Final Portfolio Value (\$mm)	Annualized Return (%)	Sharpe Ratio
ApproxEquiv	ST	1.489±0.16	12.0±3.4	0.63±0.1
	T	1.428±0.04	10.6±3.8	0.60±0.1
ExactEquiv	ST	1.411±0.15	10.3±3.4	0.62±0.2
	T	1.307±0.18	7.8±4.3	0.50±0.3
NonEquiv		1.378±0.05	9.6±1.3	0.62±0.1
Uniform		1.412	10.4	0.71
^DJI		1.293	7.7	0.53

Dow Jones index from 2001-01-01 to 2024-07-01 (see Appendix D.2 for sample data). We split the train, validation, and test data into time periods 2001-01-01-2019-01-01, 2019-01-01 - 2021-01-01, and 2021-01-01-2024-07-01, respectively. Unlike Liu et al. (2018), who used only the current timestep, we use a sliding window approach and use the previous 9 timesteps for the state. See Appendix D.2 for a more detailed description.

**Models** For this domain, we use SAC (Haarnoja et al., 2018) as our RL algorithm and consider equivariance to both the translation group and scale-translation group across the time dimension. Temporal translations can be useful as the most recent history of stock prices inform your actions, and this information may be approximately preserved across time. Temporal scaling could also be beneficial as there could be market seasonality, which is only approximately shared across different time scales. As our actions do not affect stock prices, which in turn is directly correlated with the reward, we learn an approximately invariant policy and invariant critic for both symmetry groups. As before we compare approximately equivariant, strictly equivariant, and unconstrained models. We evaluate each method on the final portfolio value (equivalent to the total episode reward), annualized return, and the Sharpe ratio (Sharpe, 1994), which is a standard financial metric that measures an asset’s risk-adjusted performance. We also include as baselines a uniform holding strategy **Uniform**, where we initially buy equal values of each stock and hold, and the Dow Jones index **^DJI**.

**Results** Table 2 lists the average test results of the learned policies on the stock trading domain. The **ApproxEquiv** model for both translation (T) and scale-translation (ST) outperform all baselines, with annualized returns of 10.6% and 12.0% respectively. The **Exact ST-Equiv** model outperforms **NonEquiv**, while



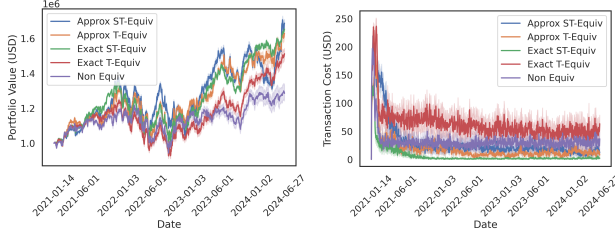


Figure 6: 10 episode rollouts from the best performing policy for each method. Approx ST-Equiv often achieves the highest portfolio value for each time step and incurs minimal transaction costs.

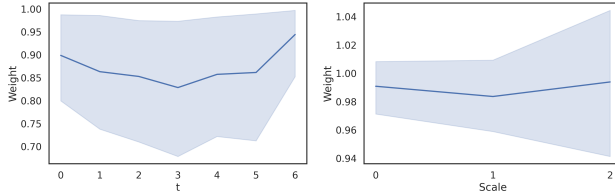


Figure 7: Visualization of relaxed weights for translation (left) and scale (right) over all runs. The relaxed weights for translation differ for each timestep and are similar for scale.

the **Exact T-Equiv** model does worse. These observations suggest that temporal scale and translation symmetries can be good biases in analyzing financial data and that translation symmetry may be more approximate than scale. We also visualize 10 episode rollouts of the best-performing policies in Figure 6, with the portfolio values on the left and transaction costs on the right. The Approx ST-Equiv method achieves the highest portfolio value for most timesteps and incurs lower transaction costs than the exactly equivariant policies. We note that overall the annualized returns are fairly low, as the test dataset from 2021-01-01 to 2024-07-01 includes both the COVID-19 pandemic and 2022 stock market decline.

We visualize the relaxed weights of the first layer of the encoder across translation (left) and scale (right) in Figure 7. For translation, our model places higher weights on the very last timestep. This matches our intuition as the most recent stock prices and portfolio holdings would be most informative in determining the optimal action. For scale, we find that the relaxed weights do not differ greatly, but there is increased variance with increasing scale.

## 6 DISCUSSION

We proposed a novel approximately equivariant architecture using relaxed group convolutions for model-free reinforcement learning. Our experimental results on continuous control domains and a stock trading

problem with real-world data demonstrate that the approximately equivariant model performs similarly to an exactly equivariant model in domains with perfect symmetry but outperforms it in most domains with symmetry breaking factors. This suggests that our method can act as a much more flexible alternative to exactly equivariant agents that can boost sample efficiency in a wider variety of settings and is also more robust to perturbations.

**Limitations and Future Work** While we did consider real-world data in the stock trading domain, our continuous control domains used simplified observations and synthetic symmetry breaking. Furthermore, exactly equivariant networks perform better in some modified domains than others (**Reacher** vs. **Acrobot**). Another limitation is that, as with all equivariant networks, the symmetry group and how it acts on the state and action spaces need to be known in advance. An interesting future direction could be to quantify exactly what types of symmetry breaking factors could lead to higher performance for approximately equivariant RL, possibly by measuring equivariance error. Other future work includes proving bounds on the optimal policy  $\pi(s)$  and  $\pi(gs)$  or applying approximately equivariant RL in robotic manipulation, where kinematic constraints or obstacles can break symmetry.

## Acknowledgments

This project was supported in part by NSF grants #2134178, 2107256, 2314182. The authors thank Hyunwoo Ryu for helpful discussions and pointing us to the Residual Pathways Prior baseline.

## Disclaimer

This paper was prepared for informational purposes [“in part” if the work is collaborative with external partners] by the Artificial Intelligence Research group of JPMorgan Chase & Co. and its affiliates (“JP Morgan”) and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

## References

- “Yahoo! Finance”. <https://finance.yahoo.com/>, 1997. Accessed: 2024-07-31.
- Fischer Black. Noise. *The journal of finance*, 41(3): 528–543, 1986.
- Alexander Bogatskiy, Brandon Anderson, Jan Offermann, Marwah Roussi, David Miller, and Risi Kondor. Lorentz group equivariant neural network for particle physics. In *International Conference on Machine Learning*, pages 992–1002. PMLR, 2020.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR, 2016.
- Taco S Cohen and Max Welling. Steerable cnns. In *International Conference on Learning Representations*, 2017.
- Luca Falorsi, Pim De Haan, Tim R Davidson, Nicola De Cao, Maurice Weiler, Patrick Forré, and Taco S Cohen. Explorations in homeomorphic variational auto-encoding. *arXiv preprint arXiv:1807.04689*, 2018.
- Marc Finzi, Gregory Benton, and Andrew G Wilson. Residual pathway priors for soft equivariance constraints. *Advances in Neural Information Processing Systems*, 34:30037–30049, 2021.
- Fabian Fuchs, Daniel Worrall, Volker Fischer, and Max Welling. Se (3)-transformers: 3d roto-translation equivariant attention networks. *Advances in neural information processing systems*, 33:1970–1981, 2020.
- Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp: Learning continuous latent space models for representation learning. In *International conference on machine learning*, pages 2170–2179. PMLR, 2019.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- Lingshen He, Yuxuan Chen, Zhengyang Shen, Yibo Yang, and Zhouchen Lin. Neural epdos: Spatially adaptive equivariant partial differential operator based networks. In *The Eleventh International Conference on Learning Representations*, 2022.
- Hyunsu Kim, Hyungi Lee, Hongseok Yang, and Juho Lee. Regularizing towards soft equivariance under mixed symmetries. In *International Conference on Machine Learning*, pages 16712–16727. PMLR, 2023.
- David M Knigge, David W Romero, and Erik J Bekkers. Exploiting redundancy: Separable group convolutional networks on lie groups. In *International Conference on Machine Learning*, pages 11359–11386. PMLR, 2022.
- Xiao-Yang Liu, Zhuoran Xiong, Shan Zhong, Hongyang Yang, and Anwar Walid. Practical deep reinforcement learning approach for stock trading. *arXiv preprint arXiv:1811.07522*, 2018.
- Daniel McNeela. Almost equivariance via lie algebra convolutions. In *NeurIPS 2023 Workshop on Symmetry and Geometry in Neural Representations*, 2023.
- Arnab Kumar Mondal, Pratheeksha Nair, and Kaleem Siddiqi. Group equivariant deep reinforcement learning. *arXiv preprint arXiv:2007.03437*, 2020.
- Alfred Müller. How does the value function of a markov decision process depend on the transition probabilities? *Mathematics of Operations Research*, 22(4):872–885, 1997.
- Jung Yeon Park, Ondrej Biza, Linfeng Zhao, Jan-Willem Van De Meent, and Robin Walters. Learning symmetric embeddings for equivariant world models. In *International Conference on Machine Learning*, pages 17372–17389. PMLR, 2022.
- Mircea Petrache and Shubhendu Trivedi. Approximation-generalization trade-offs under (approximate) group equivariance. *Advances in Neural Information Processing Systems*, 36, 2024.
- Robin Quessard, Thomas Barrett, and William Clements. Learning disentangled representations and group structure of dynamical environments. *Advances in Neural Information Processing Systems*, 33:19727–19737, 2020.
- Balaraman Ravindran and Andrew G Barto. Model minimization in hierarchical reinforcement learning. In *Abstraction, Reformulation, and Approximation: 5th International Symposium, SARA 2002 Kananaskis, Alberta, Canada August 2–4, 2002 Proceedings 5*, pages 196–211. Springer, 2002.
- David W Romero and Suhas Lohit. Learning partial equivariances from data. *Advances in Neural Information Processing Systems*, 35:36466–36478, 2022.
- Ashwin Samudre, Mircea Petrache, Brian D Nord, and Shubhendu Trivedi. Symmetry-based structured matrices for efficient approximately equivariant networks. *arXiv preprint arXiv:2409.11772*, 2024.
- Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks. In *International conference on machine learning*, pages 9323–9332. PMLR, 2021.
- William F Sharpe. The sharpe ratio. *Journal of portfolio management*, 21(1):49–58, 1994.

- Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert RG Lanckriet. On integral probability metrics,  $\phi$ -divergences and binary classification. *arXiv preprint arXiv:0901.2698*, 2009.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1802.00690*, 2018.
- Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lussann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- Tycho van der Ouderaa, David W Romero, and Mark van der Wilk. Relaxing equivariance constraints with non-stationary continuous filters. *Advances in Neural Information Processing Systems*, 35:33818–33830, 2022.
- Elise Van der Pol, Daniel Worrall, Herke van Hoof, Frans Oliehoek, and Max Welling. Mdp homomorphic networks: Group symmetries in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:4199–4210, 2020.
- Dian Wang, Jung Yeon Park, Neel Sortur, Lawson LS Wong, Robin Walters, and Robert Platt. The surprising effectiveness of equivariant models in domains with latent symmetry. In *The Eleventh International Conference on Learning Representations*, 2022a.
- Dian Wang, Robin Walters, and Robert Platt. SO(2)-equivariant reinforcement learning. In *International Conference on Learning Representations*, 2022b. URL [https://openreview.net/forum?id=7F9c0hdvfk\\_](https://openreview.net/forum?id=7F9c0hdvfk_).
- Dian Wang, Xupeng Zhu, Jung Yeon Park, Mingxi Jia, Guanang Su, Robert Platt, and Robin Walters. A general theory of correct, incorrect, and extrinsic equivariance. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Rui Wang, Robin Walters, and Rose Yu. Incorporating symmetry into deep dynamics models for improved generalization. In *International Conference on Learning Representations*, 2020.
- Rui Wang, Robin Walters, and Rose Yu. Approximately equivariant networks for imperfectly symmetric dynamics. In *International Conference on Machine Learning*, pages 23078–23091. PMLR, 2022c.
- Rui Wang, Elyssa Hofgard, Han Gao, Robin Walters, and Tess Smidt. Discovering symmetry breaking in physical systems with relaxed group convolution. In *Forty-first International Conference on Machine Learning*, 2024b.
- Maurice Weiler and Gabriele Cesa. General e (2)-equivariant steerable cnns. *Advances in neural information processing systems*, 32, 2019.
- Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco S Cohen. 3d steerable cnns: Learning rotationally equivariant features in volumetric data. *Advances in Neural Information Processing Systems*, 31, 2018.
- Matthias Weissenbacher, Rishabh Agarwal, and Yoshinobu Kawahara. Sit: Symmetry-invariant transformers for generalisation in reinforcement learning. In *International Conference on Machine Learning*, pages 52695–52719. PMLR, 2024.
- Marysia Winkels and Taco S Cohen. 3d g-cnns for pulmonary nodule detection. *arXiv preprint arXiv:1804.04656*, 2018.
- Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhambetov, and Gabriel J Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5028–5037, 2017.
- Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021.
- Xupeng Zhu, Dian Wang, Ondrej Biza, Guanang Su, Robin Walters, and Robert Platt. Sample efficient grasp learning using equivariant models. In *Robotics: Science and Systems*, 2022.
- Martin Zinkevich and Tucker R Balch. Symmetry in markov decision processes and its implications for single agent and multiagent learning. In *Proceedings of the Eighteenth International Conference on Machine Learning*, page 632, 2001.

## Checklist

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [No]
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes]
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. [Yes]
  - (b) Complete proofs of all theoretical results. [Yes]
  - (c) Clear explanations of any assumptions. [Yes]
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes]
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes]
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (a) Citations of the creator If your work uses existing assets. [Yes]
  - (b) The license information of the assets, if applicable. [Yes]
  - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes]
  - (d) Information about consent from data providers/curators. [No]
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
  - (a) The full text of instructions given to participants and screenshots. [Not Applicable]
  - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
  - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]



# Approximate Equivariance in Reinforcement Learning: Supplementary Materials

## A PROOF OF THEOREM 1

The proof is established by first deriving the deviation for an (arbitrary) finite-horizon discounted problem and then using this to derive the bounds for the infinite horizon case. All intermediate results are collected as propositions.

Consider an discounted-reward finite-horizon (horizon length is  $T$  rather than infinity) MDP with the same reward and transition kernel as the original MDP (independent of time). For a given stochastic policy  $\pi = (\pi_1, \pi_2 \dots, \pi_{T-1})$ , let

$$\begin{aligned}\mathcal{V}_t^\pi(s) &= \mathbb{E}^\pi \left[ \sum_{k=t}^{T-1} \gamma^{k-t} R(s_k, a_k) \middle| s_t = s \right], \\ \mathcal{Q}_t^\pi(s, a) &= \mathbb{E}^{s_{t+1}} \left[ R(s, a) + \gamma \mathcal{V}_{t+1}^\pi(s_{t+1}) \middle| s_t = s, a_t = a \right],\end{aligned}$$

be the finite-horizon counterparts of the expected return and action-value. Recursively define the policy independent counterparts as follows:

$$\begin{aligned}\mathcal{V}_T(s_T) &= 0, \quad \mathcal{V}_T(g s_T) = 0, \\ \mathcal{Q}_t(s_t, a_t) &= R(s_t, a_t) + \gamma \int_S \mathcal{V}_{t+1}(s_{t+1}) P(s_{t+1} | s_t, g_t), \\ \mathcal{Q}_t(g s_t, g a_t) &= R(g s_t, g a_t) + \gamma \int_S \mathcal{V}_{t+1}(g s_{t+1}) P(g s_{t+1} | g s_t, g a_t), \\ \mathcal{V}_t(s_t) &= \sup_{a_t \in \mathcal{A}} \mathcal{Q}_t(s_t, a_t), \quad \mathcal{V}_t(g s_t) = \sup_{a_t \in \mathcal{A}} \mathcal{Q}_t(g s_t, g a_t).\end{aligned}$$

We also define

$$V_t^\pi(s) := \mathbb{E}^\pi \left[ \sum_{k=t}^{\infty} \gamma^{k-t} R(s_k, a_k) \middle| s_t = s \right], \quad Q_t^\pi(s, a) := \mathbb{E}^\pi \left[ \sum_{k=t}^{\infty} \gamma^{k-t} R(s_k, a_k) \middle| s_t = s, a_t = a \right].$$

Let  $V_t(s_t) := \sup_{\pi} V_t^\pi(s_t)$  and

$$Q_t(s_t, a_t) = \mathbb{E} \left[ R(s_t, a_t) + \gamma V_{t+1}(s_{t+1}) \middle| s_t, a_t \right].$$

Note that  $V_t^\pi$  and  $Q_t^\pi$  equal  $V^\pi$  and  $Q^\pi$  defined in (5) for any  $t$ , and  $V_t$  and  $Q_t$  equal  $V^*$  and  $Q^*$ . We introduce the notations only to make the connection between the finite-horizon and infinite-horizon MDPs clearer.

**Proposition 1.** *For a  $(G, \epsilon_R, \epsilon_P)$ -invariant MDP, the following holds at any  $t$ ,*

$$|\mathcal{Q}_t(s_t, a_t) - Q_t(g s_t, g a_t)| \leq \alpha_t, \quad \text{and} \quad |\mathcal{V}_t(s_t) - \mathcal{V}_t(g s_t)| \leq \alpha_t,$$

where  $\alpha_t$  is given by the following recursion:  $\alpha_{T+1} = 0$  and

$$\alpha_t = \epsilon_R + \gamma \left\{ \rho_{\mathcal{F}}(\mathcal{V}_{t+1}) \epsilon_P + \alpha_{t+1} \right\}.$$

*Proof.* We will prove the results using induction. First, note that the result is true for  $T$  by definition. Suppose the result is true for  $t + 1$ , and consider the differential at time  $t$ ,

$$\begin{aligned}
 |\mathcal{Q}_t(s_t, a_t) - \mathcal{Q}_t(g s_t, g a_t)| &\leq |R(s_t, a_t) - R(g s_t, g a_t)| \\
 &\quad + \gamma \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(s_{t+1}) P(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(g s_{t+1}) P(g s_{t+1} | g s_t, g a_t) \right| \\
 &\leq \epsilon_R + \gamma \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(s_{t+1}) P(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(g s_{t+1}) P(s_{t+1} | s_t, a_t) \right| \\
 &\quad + \gamma \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(g s_{t+1}) P(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(g s_{t+1}) P(g s_{t+1} | g s_t, g a_t) \right| \\
 &\leq \epsilon_R + \gamma \rho_{\mathcal{P}}(\mathcal{V}_{t+1}) \epsilon_P + \gamma \int_{\mathcal{S}} \left| \mathcal{V}_{t+1}(s_{t+1}) - \mathcal{V}_{t+1}(g s_{t+1}) \right| P(s_{t+1} | s_t, a_t).
 \end{aligned}$$

The last inequality follows by using the decomposition using Minkowski's functional. Further, note that

$$\left| \mathcal{V}_{t+1}(s_{t+1}) - \mathcal{V}_{t+1}(g s_{t+1}) \right| \leq \sup_{a_{t+1} \in \mathcal{A}} |\mathcal{Q}_{t+1}(s_{t+1}, a_{t+1}) - \mathcal{Q}_{t+1}(g s_{t+1}, g a_{t+1})| \leq \alpha_{t+1},$$

by induction assumption, and the fact that when  $g$  is onto

$$\sup_{a' \in g\mathcal{A}} \mathcal{Q}_{t+1}(g s_t, a') = \sup_{a \in \mathcal{A}} \mathcal{Q}_{t+1}(g s_t, g a).$$

The result follows. □

**Proposition 2.** *Let the rewards  $R \in [R_{\min}, R_{\max}]$ . For an arbitrary, but finite, horizon  $T$*

$$\mathcal{Q}_t(s_t, a_t) + \frac{\gamma^{T-t}}{1-\gamma} R_{\min} \leq Q_t(s_t, a_t) \leq \mathcal{Q}_t(s_t, a_t) + \frac{\gamma^{T-t}}{1-\gamma} R_{\max}.$$

*Proof.* We have by definition,

$$\begin{aligned}
 Q_t(s_t, a_t) &= \mathbb{E} \left[ \sum_{k=t}^{\infty} \gamma^{k-t} R(s_k, a_k) \middle| s_t = s, a_t = a \right] \\
 &= \mathbb{E} \left[ R(s_t, a_t) + \gamma \mathbb{E} \left[ \sum_{k=t+1}^{\infty} \gamma^{k-(t+1)} R(s_k, a_k) \middle| s_{t+1} \right] \middle| s_t = s, a_t = a \right] \\
 &\leq \mathbb{E} \left[ R(s_t, a_t) + \gamma \mathbb{E} \left[ \mathcal{V}_{t+1}(s_{t+1}) + \frac{\gamma^{T-(t+1)} R_{\max}}{1-\gamma} \right] \middle| s_{t+1} \right] \middle| s_t = s, a_t = a \\
 &= \mathcal{Q}_t(s_t, a_t) + \frac{\gamma^{T-t}}{1-\gamma} R_{\max}.
 \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 Q_t(s_t, a_t) &= \mathbb{E} \left[ \sum_{k=t}^{\infty} \gamma^{k-t} R(s_k, a_k) \middle| s_t = s, a_t = a \right] \\
 &\geq \mathbb{E} \left[ \sum_{k=t}^{T-1} \gamma^{k-t} R(s_k, a_k) + \sum_{k=T}^{\infty} \gamma^{k-t} R_{\min} \middle| s_t, a_t \right] \\
 &= \mathbb{E} \left[ R(s_t, a_t) + \gamma \sum_{k=t+1}^{T-1} \gamma^{k-(t+1)} R(s_k, a_k) \middle| s_t, a_t \right] + \frac{\gamma^{T-t}}{1-\gamma} R_{\min} \\
 &= \mathcal{Q}_t(s_t, a_t) + \frac{\gamma^{T-t}}{1-\gamma} R_{\min}.
 \end{aligned}$$

□

We now prove Theorem 1. Let  $\mathcal{B}(\mathcal{S})$  denote the Banach space of bounded real-valued functions on  $\mathcal{S}$ . We define the Bellman optimality operator  $\mathcal{B} : \mathcal{B}(\mathcal{S}) \rightarrow \mathcal{B}(\mathcal{S})$  such that for any uniformly bounded function  $V \in \mathcal{B}(\mathcal{S})$ ,

$$\mathcal{B}V(s) = \sup_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \int_{\mathcal{S}} V(s') P(s'|s, a) \right\} \quad \forall s \in \mathcal{S}. \quad (6)$$

It is known that  $V^*$  is the (unique) fixed point of  $\mathcal{B}$ , i.e.,  $\mathcal{B}V^* = V^*$ . We note that  $V^*$  also satisfies the following equation for any  $s$

$$V^*(gs) = \sup_{a \in \mathcal{A}} \left\{ R(gs, ga) + \gamma \int_{\mathcal{S}} V^*(gs') P(gs'|gs, ga) \right\}. \quad (7)$$

To see this, consider the following arguments.

$$\begin{aligned} Q^*(s, a) &= R(s, a) + \gamma \sup_{a' \in \mathcal{A}} \int_{s' \in \mathcal{S}} Q^*(s, a) P(s'|s, a), \\ Q^*(gs, ga) &= R(gs, ga) + \gamma \sup_{a' \in \mathcal{A}} \int_{s' \in \mathcal{S}} Q^*(gs, ga) P(s'|gs, ga). \end{aligned}$$

Since  $g \in G$  permutes the elements of  $G$ , re-indexing the integral using  $\tilde{s}' = gs'$ , we have

$$\begin{aligned} Q^*(gs, ga) &= R(gs, ga) + \gamma \sup_{\tilde{a} \in g\mathcal{A}} \int_{\tilde{s}' \in g\mathcal{S}} Q(\tilde{s}', \tilde{a}') P(\tilde{s}'|gs, ga). \\ \therefore Q^*(gs, ga) &= R(gs, ga) + \gamma \sup_{a' \in \mathcal{A}} \int_{s' \in \mathcal{S}} Q(gs', ga') P(gs'|gs, ga). \end{aligned}$$

### Proof of Theorem 1:

Consider a sequence of value functions  $\mathcal{V}^{(n)}$  on the symmetry transformed domain as follows:  $\mathcal{V}^{(0)}(gs) = 0$  and  $\mathcal{V}^{(n+1)} = \mathcal{B}\mathcal{V}^{(n)}$ . For an arbitrary  $T$ , we have using Proposition 1 for any  $t \in \{1, \dots, T\}$ ,

$$|\mathcal{V}_t(s_t) - \mathcal{V}_t^{(T-t)}(gs_t)| \leq \alpha_t,$$

where

$$\alpha_t = \epsilon_R + \sum_{\tau=t+1}^{T-1} \gamma^{\tau-t} [\rho_{\mathcal{F}}(\mathcal{V}^{(T-\tau)}) \epsilon_P + \epsilon_R].$$

From Proposition 2, we have, noting that  $\mathcal{V}(s) = \sup_a Q(s, a)$ , that

$$\mathcal{V}_t^{(T-t)}(gs_t) - \alpha_t + \frac{\gamma^{T-t}}{1-\gamma} R_{\min} \leq V_t(s_t) \leq \mathcal{V}_t^{(T-t)}(gs_t) + \alpha_t + \frac{\gamma^{T-t}}{1-\gamma} R_{\max}$$

By Banach fixed point theorem, we know that  $\lim_{T \rightarrow \infty} \mathcal{V}_t^{(T-t)} = V^*$ . By continuity of  $\rho_{\mathcal{F}}(\cdot)$ , we have that  $\lim_{T \rightarrow \infty} \rho_{\mathcal{F}}(\mathcal{V}^{(T-\tau)}) = \rho_{\mathcal{F}}(V^*)$  whence  $\lim_{T \rightarrow \infty} \alpha_t = \alpha := \frac{\epsilon_R + \gamma \rho_{\mathcal{F}}(V^*) \epsilon_P}{1-\gamma}$ . Therefore, taking the limit, we have

$$V^*(gs_t) - \alpha \leq V_t(s_t) \leq V^*(gs_t) + \alpha.$$

A similar argument establishes the result for  $Q$  using the onto function  $g$ . The claims in Theorem 1 follows by recognizing that  $V_t$  and  $Q_t$  exactly equal  $V^*$  and  $Q^*$ . □

## B CASE OF FINITE HORIZON: NO DISCOUNTING

As is clear from Theorem 1, when  $\gamma \rightarrow 1$ , the bound becomes trivial and not useful. In this section we will briefly discuss the case when the discount factor  $\gamma = 1$ . In this setting, we allow the transition functions to be a function of  $t$ .

**Proposition 3.** Let  $|R(gs_t, ga_t) - R(s_t, a_t)| \leq \epsilon_R$  and  $d_{\mathcal{F}}(P_t(gs'_t | gs_t, ga_t), P_t(s'_t | s_t, a_t)) \leq \epsilon_P(t)$ . For a finite-horizon MDP of duration  $T$ , we have

$$|\mathcal{Q}_t(s_t, a_t) - \mathcal{Q}_t(gs_t, ga_t)| \leq \alpha_t, \quad |\mathcal{V}_t(s_t) - \mathcal{V}_t(gs_t)| \leq \alpha_t$$

where  $\alpha_{T+1} = 0$  and for  $t \in \{1, 2, \dots, T\}$ ,

$$\alpha_t = \epsilon_R + \sum_{\tau=t+1}^T \left[ \rho_{\mathcal{F}}(\mathcal{V}_{\tau}) \epsilon_P(\tau - 1) + \epsilon_R \right].$$

*Proof.* The proof proceeds as in Proposition 1. We have

$$\begin{aligned} |\mathcal{Q}_t(s_t, a_t) - \mathcal{Q}_t(gs_t, ga_t)| &\leq |R(s_t, a_t) - R(gs_t, ga_t)| \\ &\quad + \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(s_{t+1}) P_t(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(gs_{t+1}) P_t(gs_{t+1} | gs_t, ga_t) \right| \\ &\leq \epsilon_R + \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(s_{t+1}) P_t(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(gs_{t+1}) P_t(s_{t+1} | s_t, a_t) \right| \\ &\quad + \left| \int_{\mathcal{S}} \mathcal{V}_{t+1}(gs_{t+1}) P_t(s_{t+1} | s_t, a_t) - \int_{\mathcal{S}} \mathcal{V}_{t+1}(gs_{t+1}) P_t(gs_{t+1} | gs_t, ga_t) \right| \\ &\leq \epsilon_R + \rho_{\mathcal{F}}(\mathcal{V}_{t+1}) \epsilon_P(t) + \alpha_{t+1} := \alpha_t. \end{aligned}$$

The result follows by recursion.  $\square$

## C BACKGROUND AND METHOD

### C.1 Equivariance with Group Convolutions

Group convolutions (Cohen and Welling, 2016) generalize standard convolutions, which are translation-equivariant, to be equivariant to a group  $G$ . Group convolutions act on signals over the group  $f : G \rightarrow \mathbb{R}$ . As many data samples are not natively of this form (e.g. an image), the input must first be lifted onto a function in  $G$ . For example, let  $f_0 : \mathbb{Z}^2 \rightarrow \mathbb{R}$  be the input signal, a grayscale image, and  $H = D_2$  be the group. The lifting convolution lifts  $f_0$  from  $\mathbb{Z}^2$  to  $G = D_2 \times \mathbb{Z}^2$  by

$$(f_0 \star \psi)(x, h) = \sum_{y \in \mathbb{Z}^2} f_0(y) \psi(h^{-1}(y - x)), \quad (8)$$

where  $h \in H$ . Practically, the lift operation creates  $|H|$ , the order of group  $H$ , images by acting on  $x$  by  $h^{-1}$ . Typically the lift operation is the first layer of the network, followed by subsequent group convolutions, nonlinearities, or other equivariant layers. We use relaxed versions of the lift and group convolutions as described in Wang et al. (2022c) and the main paper.

### C.2 Steerable Convolutions

As an alternative to group convolutions, one can use steerable convolutions (Weiler et al., 2018) that use weight tying to generalize to continuous groups and are more parameter-efficient. Let  $H < O(2)$  be the subgroup which acts on  $\mathbb{R}^2$  by matrix multiplication on the input and output channel spaces  $\mathbb{R}^c$  and  $\mathbb{R}^d$  by  $\rho_{\text{in}}$  and  $\rho_{\text{out}}$ , respectively. Then  $G = H \ltimes \mathbb{R}^2$ . Given input signal  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^c$ , then standard convolution over  $\mathbb{R}^2$  with kernel  $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}^{d \times c}$  is  $G$ -equivariant if  $\psi$  satisfies

$$\psi(hx) = \rho_{\text{out}}(g) \psi(x) \rho_{\text{in}}(h^{-1}), \quad (9)$$

for all  $h \in H$ . Intuitively, this kernel constraint ensures that the output features transform by  $\rho_{\text{out}}$  when the input features are transformed by  $\rho_{\text{in}}$ . Kernels that satisfy this constraint have been solved for many common subgroups of  $E(2)$ , see Weiler and Cesa (2019) for more details.

Using the example of grayscale images as in Section C.1, let the input feature be  $f : \mathbb{Z}^2 \rightarrow \mathbb{R}$  and  $\{\psi_k\}_{k=1}^K$  be a precomputed, nontrainable equivariant kernel basis of  $K$  kernels that satisfy Eq. (9). Assume that both the



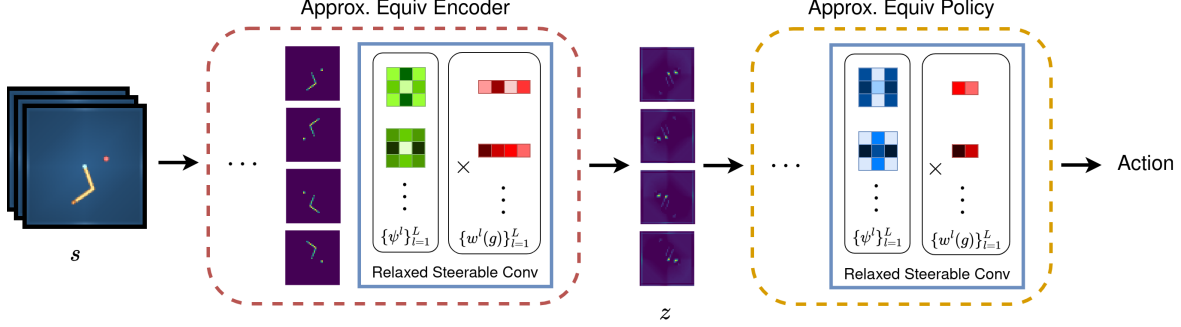


Figure 8: Illustration of an approximately  $D_2$ -equivariant encoder and policy using relaxed steerable convolution layers. The critic is not shown and is approximately invariant.

number of input and output channels is 1 and let  $w \in \mathbb{R}^K$  be the trainable coefficients of the kernels. Then a  $G$ -steerable convolution is defined as

$$(f \star \psi)(x) = \sum_{y \in \mathbb{Z}^2} \sum_{k=1}^K (w_k \psi_k(y)) f(x + y), \quad (10)$$

where  $x \in \mathbb{Z}^2$  is the spatial position and  $w_k$  is the weight associated with kernel  $\psi_k$ .

**Relaxed Steerable Convolution** As described in Wang et al. (2022c), one can also use relaxed versions of steerable convolutions by letting the trainable weights  $w$  also depend on  $y$ . A relaxed  $G$ -steerable convolution is defined as

$$(f \tilde{\star} \psi)(x) = \sum_{y \in \mathbb{Z}^2} \sum_{k=1}^K (w_k(y) \psi_k(y)) f(x + y). \quad (11)$$

Allowing the trainable weights  $w_k$  to also depend on the absolute spatial position  $y$  breaks the equivariance constraint in Eq. (9).

By replacing relaxed group convolutions with relaxed steerable convolutions, we can also design a variant of our proposed approximately equivariant RL architecture (Figure 8).

## D EXPERIMENT DETAILS

### D.1 Continuous Control

**Acrobot** We use the **swingup** task. The domain consists of two joints where the goal is to apply torque to the inner joint so that both joints are near vertical. We use  $D_1$  as the symmetry group, i.e. vertical reflection, and the action  $a \in \mathbb{R}$  transforms via the sign representation  $\rho_{\text{sign}}$ , where  $\rho_{\text{sign}}(\text{flip})(a) = -a$ . For variants, we consider 1) **repeat\_action** - the action is repeated when the inner joint is in the fourth quadrant and 2) **gravity** - gravity  $\vec{g} = [0, 0, -9.81]$  is modified to  $[-2, 2, -9.81]$ .

**CartPole** We consider the **swingup** task. The domain consists of a pole swinging on a cart and the goal is to move the cart left or right ( $a \in \mathbb{R}$ ) to make the pole upright. The symmetry group and action representation are the same as in **Acrobot**,  $D_1$  and  $\rho_{\text{sign}}$ . For variants, we consider 1) **repeat\_action** - the action is repeated when the pole is in the first quadrant, 2) **gravity** - gravity is modified to  $[0.2, -0.2, -9.81]$ , and 3) **reflect\_action** - the pole angle is in  $[0, \frac{\pi}{4}]$ . Gravity is modified less than in **Acrobot** as too high values forced the cart out of frame.

**Cup Catch** The domain consists of a ball attached to the bottom of the cup and the goal is to move the cup to catch the ball inside the cup. The action  $(x, z) \in \mathbb{R}^2$  is the cup’s spatial position. The symmetry group is  $D_1$  and the action representation is  $\rho_{\text{sign}} \oplus \rho_0$ , where the  $x$  position transforms via the sign representation and the

$z$  transforms via the trivial representation  $\rho_0$ . For variants, we consider 1) **repeat\_action** - the ball  $x$  position greater than 0.0 and  $z$  position is greater than 0.3, 2) **gravity** - gravity is modified to  $[-2, 2, -9.81]$ , and 3) **reflect\_action** - same as **repeat\_action**.

**Reacher** We consider the **hard** task. The domain consists of two joints and the goal is to apply torques to make the end effector reach the target. The action  $a \in \mathbb{R}^2$ . The symmetry group is  $D_2$ , i.e. vertical reflections and  $\pi$  rotations, and the action transforms via the quotient representation  $2\rho_{\text{quot}}$ , where the torques for both joints are invariant to rotations and flip signs for vertical reflections. For variants, we consider 1) **repeat\_action** - the inner joint angle is in  $[0, \frac{\pi}{2}]$  and 2) **reflect\_action** - the inner joint angle is in  $[\frac{\pi}{2}, \pi]$ .

### D.1.1 Training Details

For all DeepMind Control Suite (DMC) domains, we fix the episode length to 1000 and use RGB image of size  $85 \times 85$ . We considered four domains of varying difficulty, of which **Acrobot** is the hardest. In the original DrQv2 implementation (Yarats et al., 2021), the encoder reduces the spatial dimensions to  $35 \times 35$ , which is then flattened to be input to the policy and critic. We follow Wang et al. (2022a) and further reduce the spatial dimensions to  $7 \times 7$  for faster training for all models. We reduce the replay buffer size from 1,000,000 to 500,000 to slightly reduce the memory footprint. All other hyperparameters are kept the same as in Yarats et al. (2021).

For the exactly equivariant and approximately equivariant models, we reduce the number of channels by  $\sqrt{|G|}$  where  $|G|$  is the order of the group to preserve roughly the same number of parameters as the non-equivariant model. We use  $L = 1$  filters for the approximately equivariant model in all experiments.

RPP contains both the non-equivariant layers and exactly equivariant layers and thus has roughly twice as many parameters as **ExactEquiv**. For the critic moving average speed  $\tau$ , we use the default  $\tau = 0.01$  for **CartPole** and **Reacher** and  $\tau = 0.009$  for **Acrobot** and **Ball in Cup**.

The plots in Figure 3 show the mean reward of 10 episodes, evaluated every 20,000 environment steps. For the results in Table 1, we use  $\sigma = 0.02$  for **Acrobot** and **Reacher** and  $\sigma = 0.06$  for **CartPole** and **Ball in Cup**.

The continuous control experiments were run on single GPUs of different types. **Acrobot** was run on an Nvidia RTX 4090 and all other experiments were run on an Nvidia RTX 2080 Ti. We note that the wall clock time for training both exactly and approximately equivariant models is longer than that for a non equivariant model, even though they are generally more sample efficient. This is because equivariant neural networks often incur more overhead in implementation - for group convolutions, the kernel must be transformed and the outputs must be stacked and for steerable convolutions, the basis must be projected onto matrices at every forward pass.

## D.2 Stock Trading

We formulate the stock trading problem as an MDP as described in Liu et al. (2018). The state consists of the cash balance  $c_t$ , the stock prices  $p_t^n$ , the number of shares in the current portfolio  $h_t^n$ , and other technical indicators  $i_t^n$  for time  $t$  stock  $n \in \{1, \dots, N\}$ . The actions  $x_t^n$  are the number of stocks to buy and sell for each stock  $n$  and are bounded to  $[-M, M]$  where  $M$  was set to 100. The reward  $r_t$  is the scaled difference in portfolio values between consecutive timesteps and we assume that the market dynamics are not affected by our trading. There is a small transaction cost  $\epsilon^n = 0.001$  for every trade. Initially, the portfolio contains 0 shares and the cash balance is 1,000,000. This can be formulated as a constrained program as follows

$$\begin{aligned}
 \max \quad & \sum_t r_t \\
 \text{s.t.} \quad & -M \leq a_t^n \leq M, & \forall n, t \\
 & a_t^n \geq -h_t^n, & \forall n, t \\
 & a_t^n \leq \lfloor c_t / (p_t^n (2 + \epsilon^n)) \rfloor & \forall n, t \\
 & c_t \geq 0 & \forall t \\
 & c_{t+1} = c_t - \sum_n a_t^n p_t^n (1 + \epsilon^n) & \forall t \\
 & h_{t+1}^n = h_t^n + a_t^n & \forall n, t
 \end{aligned}$$

$$\begin{aligned}
 r_{t+1} &= (c_{t+1} - c_t) + \sum_n (p_{t+1}^n h_{t+1}^n - p_t^n h_t^n) & \forall t \\
 c_0 &= 1,000,000 \\
 h_0^n &= 0 & \forall n \\
 h_t^n &\in \mathbb{Z}^+, a_t^n \in \mathbb{Z}, c_t \in \mathbb{R}^+.
 \end{aligned}$$

The financial data was pulled from Yahoo Finance (yfi, 1997) for the time period 2001-01-01 to 2024-07-01 (see Figure 9 for sample stock prices). As historical stock prices and portfolio can be important for determining the action, we use the previous  $H = 9$  timesteps for the state, unlike Liu et al. (2018).

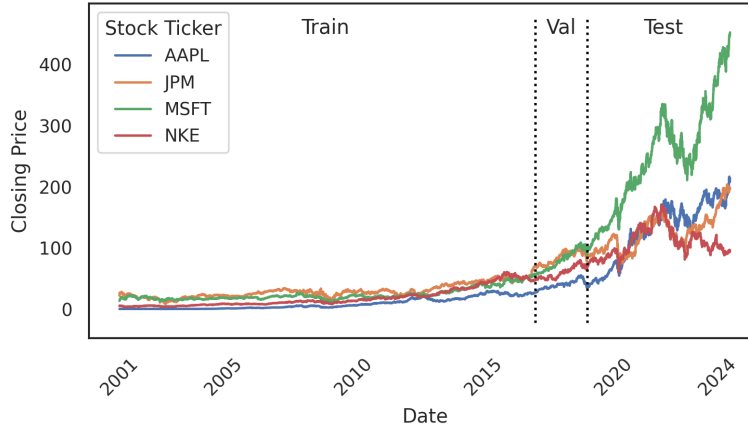


Figure 9: Sample stock trading data. We use a sliding window of the stock prices, current portfolio, cash balance, and other indicators as the state. The dataset is split into train/val/test as shown.

### D.2.1 Training Details

or all models, we use 4 layers for the shared encoder, 1 layer for the actor, and 2 layers for the critic. The non equivariant model uses linear layers after flattening the input, while the exactly equivariant and approximately equivariant models use group convolutions and relaxed group convolutions with a kernel size of 5, respectively. We consider both temporal translations and temporal scale-translations. For scale-translation, we use separable group convolutions (Knigge et al., 2022) and use 3 scale factors 0.8, 0.98, 1.2. We control the number of channels so that the total number of parameters is roughly equal to the non equivariant model. We use  $L = 1$  filters for the approximately equivariant model in all experiments.

The stock trading experiments were run on a single Nvidia RTX 2080 Ti. All other hyperparameters are given in Table 3.

Table 3: Hyperparameters used for stock trading experiments

Hyperparameter	ApproxEquiv	ExactEquiv	NonEquiv
Batch size		64	
Learning rate		1e-4	
$\alpha$		0.05	
$\tau$		0.005	
Discount factor		0.99	
Hidden dim/channels	64	64	128
Encoder output dim/channels		256	

Table 4: Total episode reward on 50 rollouts for the best policy when trained on noisy inputs and tested in the noisy domain. Gray values indicate 95% CI. **ApproxEquiv** learns a more robust policy than baselines on the modified **BallInCup** domain.

		ApproxEquiv	ExactEquiv	NonEquiv
BALL IN CUP (Noisy)	Original	971 $\pm$ 1.7	<b>977</b> $\pm$ 3.8	914 $\pm$ 7.9
	Gravity	<b>973</b> $\pm$ 2.7	952 $\pm$ 5.7	942 $\pm$ 10.

## E ROBUSTNESS WITH NOISE AUGMENTATION

Table 4 shows the results from training policies with noisy inputs and evaluating their robustness to noise at test time. This experiment tests whether our approximately equivariant method is truly more robust to noise than other baselines trained with noise augmentation. We find that our approximately equivariant method is more robust than baselines for the modified domain, even when trained with noise augmentation.