

# Leveraging symmetries in pick and place

Haojie Huang , Dian Wang, Arsh Tangri, Robin Walters\* and Robert Platt\*

The International Journal of  
Robotics Research  
2024, Vol. 43(4) 550–571  
© The Author(s) 2024  
Article reuse guidelines:  
[sagepub.com/journals-permissions](https://sagepub.com/journals-permissions)  
DOI: 10.1177/02783649231225775  
[journals.sagepub.com/home/ijr](https://journals.sagepub.com/home/ijr)



## Abstract

Robotic pick and place tasks are symmetric under translations and rotations of both the object to be picked and the desired place pose. For example, if the pick object is rotated or translated, then the optimal pick action should also rotate or translate. The same is true for the place pose; if the desired place pose changes, then the place action should also transform accordingly. A recently proposed pick and place framework known as *Transporter Net* (Zeng, Florence, Tompson, Welker, Chien, Attarian, Armstrong, Krasin, Duong, Sindhvani et al., 2021) captures some of these symmetries, but not all. This paper analytically studies the symmetries present in planar robotic pick and place and proposes a method of incorporating equivariant neural models into *Transporter Net* in a way that captures all symmetries. The new model, which we call *Equivariant Transporter Net*, is equivariant to both pick and place symmetries and can immediately generalize pick and place knowledge to different pick and place poses. We evaluate the new model empirically and show that it is much more sample-efficient than the non-symmetric version, resulting in a system that can imitate demonstrated pick and place behavior using very few human demonstrations on a variety of imitation learning tasks.

## Keywords

Deep learning, manipulation, vision

Received 13 February 2023; Revised 6 November 2023; Accepted 17 December 2023

Senior Editor: Jose-Luis Blanco

Associate Editor: Shoudong Huang

## 1. Introduction

Pick and place is an important paradigm in robotic manipulation where a complex manipulation problem can be decomposed into a sequence of grasp (pick) and place operations. Recently, multiple learning approaches have been proposed to solve this problem, including Zeng et al. (2021); Wang et al. (2021). These methods focus on a simple version of the planar pick and place problem where the method looks at the scene and outputs a single pick and a single place pose. This problem has an important structure in the form of symmetries in  $SE(2)$  that can be expressed with respect to the pick and place pose. The pick symmetry is easiest to see. If the object to be grasped is rotated (in the plane), then the optimal grasp pose clearly must also rotate. A similar symmetry exists in place pose. If an object is to be placed into an environment in a particular way, then if the environment rotates, the desired place pose must also rotate. Leveraging symmetries of the task could result in significant gains in sample efficiency (Zhu et al., 2022; Jia et al., 2023). Why is sample efficiency important in robot learning? Although robotic simulators could provide a huge amount of data that could be used to train a policy, there is an inevitable sim-to-real gap in applying the learned policy directly to real robots. On the other side, real-world robot data is expensive to

collect, and sample efficiency is crucial to learning a policy with a limited number of human demonstrations.

If we are to design a robotic learning system for pick and place, it should ideally encode the symmetries described above. This is a structure that exists in the problem and there is a possibility to simplify learning by encoding this structure into our learned solutions. The question is how to accomplish this. This paper examines the symmetries that exist in the pick and place problem by identifying invariant and equivariant equations that we would expect to be preserved. Then, we consider existing pick and place models and find that those architectures only express some but not all problem symmetries. Finally, we propose a novel pick and place model that we call *Equivariant Transporter Net* that encodes all

---

Khoury College of Computer Science, Northeastern University, Boston, MA, USA

\*Equal Advising

### Corresponding author:

Haojie Huang, Khoury College of Computer Science, Northeastern University, 360 Huntington Ave, Boston, MA 02115, USA.

Email: [huang.haoj@northeastern.edu](mailto:huang.haoj@northeastern.edu)

symmetries and shows that it outperforms models that do not preserve the relevant symmetries.

### 1.1. Symmetries in transporter net

This paper builds on top of the *Transporter Net* model (Zeng, Florence, Tompson, Welker, Chien, Attarian, Armstrong, Krasin, Duong, Sindhvani et al., 2021). Transporter Net is a sample-efficient model for learning planar pick and place behaviors through imitation learning. Compared to many other approaches (Qureshi et al., 2021; Curtis et al., 2022), it does not need to be pre-trained on the involved objects—it only needs to be trained on the given demonstrations. Transporter Net achieves sample efficiency in this setting by encoding the symmetry of the picked object into the model. Once the model learns to pick and place an object presented in one orientation, that knowledge immediately generalizes to a finite set of other pick poses. This is illustrated in Figure 1a. The left side of Figure 1(a) shows a pick-place problem where the robot must pick the orange object and place it inside the green outline. Because the model encodes the symmetry of the picked object, the ability to solve the place task on the left side of Figure 1(a) immediately implies an ability to solve the place task on the right side of Figure 1(a) where the object to be picked has been rotated. We will refer to this as a  $SO(2)$ -place symmetry. Since *Transporter Net* used a set of discrete rotations, it actually achieves  $C_n$ -place symmetry where  $C_n$  is the finite cyclic subgroup of  $SO(2)$  that contains a set of  $n$  rotations.

### 1.2. Equivariant Transporter Net

This paper analyzes the symmetries present in the pick and place problem and expands Transporter Net in the following ways. First, we constrain the pick model to be equivariant (an expression of symmetry) with respect to the  $SO(2)$  group by incorporating equivariant convolutional layers into the pick model. This is, if there is a rotation on the object to be picked, the pick pose will also rotate. We refer to this as a  $SO(2)$ -pick symmetry. The second way we extend Transporter Net is by making it equivariant with respect to changes in place orientation. That is, if the place model learns how to place an object in one orientation, that knowledge generalizes immediately to different place orientations. Our resulting placing model is equivariant both to changes in pick and place orientation, and can be viewed as a direct product of two groups,

$SO(2) \times SO(2)$  as illustrated in Figure 1(b). This expanded symmetry improves the sample efficiency of our model by enabling it to generalize over a larger set of problems. Finally, we also propose a goal-conditioned version of Equivariant Transporter Net where the desired place pose is provided to the system in the form of an image as shown in Figure 8.

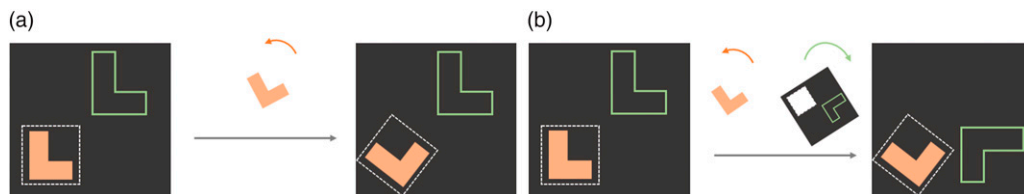
### 1.3. Contributions

Our specific contributions are as follows. (1) We systematically analyze the symmetries present in the planar pick and place problem. (2) We propose Equivariant Transporter Net, a novel version of Transporter Net that has  $C_n$ -equivariant pick symmetry and  $C_n \times C_n$ -equivariant place symmetry.<sup>1</sup> (3) We propose a variation of Equivariant Transporter Net that can be used with standard grippers rather than just suction cups. (4) We propose a goal-conditioned version of Equivariant Transporter Net. (5) We evaluate the approach both in simulation tasks and on physical robot versions of three of the gripper tasks. Our results indicate that our approach is more sample-efficient than the baselines and therefore learns better policies from a small number of demonstrations. Video and code are available at [https://haojhuang.github.io/etp\\_page/](https://haojhuang.github.io/etp_page/).

This paper extends the recent work (Huang, Wang, Walters and Platt, 2022a) in the following ways. First, we cover the concepts, algorithms, and results in a more comprehensive way. Second, we generalize our proofs of equivariance from  $C_n$  to any subgroup of  $SO(2)$ . We also analyze the extension to  $SO(3)$  mathematically and provide intuition. Third, we propose a goal-conditioned extension of the work and show that the new method outperforms on the benchmark of goal-conditioned tasks. Finally, we add an ablation study that characterizes the model for differently sized cyclic groups,  $C_n$ .

### 1.4. Comparison to related works

**1.4.1. Pick and place.** Pick and place is an important topic in manipulation. Many fundamental skills like packing, kitting, and stacking require inferring both the pick and the place action. Traditional assembly methods in factories use customized workstations so that fixed pick and place actions can be manually predefined. Recently, considerable research has focused on vision-based manipulation. Some work (Narayanan and Likhachev, 2016; Chen et al., 2019;



**Figure 1.** Visual explanation of  $SO(2)$ -equivariance (left figure) v.s.  $SO(2) \times SO(2)$ -equivariance of the place model (right figure). (a) If Transporter Network (Zeng et al., 2021) learns to place an object when it is presented in one orientation, the model is immediately able to generalize to new object orientations. (b) Our proposed Equivariant Transporter Network is able to generalize over both pick and place orientation. We view this as  $SO(2) \times SO(2)$ -place symmetry of the model.

Gualtieri and Platt, 2021) assumes that object mesh models are available in order to run ICP (Besl and McKay, 1992) and align the object model with segmented observations or completions (Yuan et al., 2018; Huang et al., 2021). Other work learns a category-level pose estimator (Yoon et al., 2003; Deng et al., 2020) or key-point detector (Nagabandi et al., 2020; Liu et al., 2020; Manuelli et al., 2019) from training on a large dataset. Recently, Wen, Lian, Bekris and Schaal (2022) realizes a close-loop intra-category policy by mimicking the extracted pose trajectory from a few video demonstrations. However, these methods often require expensive object-specific labels or pre-training, making them difficult to use widely. Recent advances in deep learning have provided other ways to rearrange objects from perceptual data. Qureshi et al. (2021) represent the scene as a graph over segmented objects to do goal-conditioned planning; Curtis et al. (2022) propose a general system consisting of a perception module, grasp module, and robot control module to solve multi-step manipulation tasks. These approaches often require prior knowledge like good segmentation module and human-level hierarchy. End-to-end models (Zakka et al., 2020; Khansari et al., 2020; Devin et al., 2020; Berscheid et al., 2020) that directly map input observations to actions can learn quickly and generalize well. Shridhar, Manuelli and Fox (2022a) learn one multi-task policy with language-conditioned imitation learning. Shridhar, Manuelli and Fox (2022b) directly extend this idea to 3D keyframe multi-task policy learning with Perceiver IO transformer (Jaegle et al., 2021). Wu et al. (2020) achieve fast learning speed on deformable-object manipulation tasks with reinforcement learning. However, most methods need to be trained on large datasets. For example, Khansari et al. (2020) collects a dataset with 7.2 million samples. Devin et al. (2020) collects 40K grasps and places per task. Zakka et al. (2020) collects 500 disassembly sequences for each kit. The focus of this paper is on improving the sample efficiency of this class of methods on various manipulation tasks.

*1.4.2. Equivariance learning in manipulation.* Fully Convolutional Networks (FCNs) are translationally equivariant and have been shown to improve learning efficiency in many manipulation tasks (Zeng, Song, Yu, Donlon, Hogan, Bauza, Ma, Taylor, Liu, Romo et al., 2018b; Morrison et al., 2018). The idea of encoding SE(2) symmetries in the structure of neural networks is first introduced in G-Convolution (Cohen and Welling, 2016). The extension work proposes an alternative architecture, Steerable CNN (Cohen and Welling, 2017). Weiler and Cesa (2019) propose a general framework for implementing E(2)-Steerable CNNs. Weiler, Geiger, Welling, Boomsma and Cohen (2018) first investigated the SE(3) steerable convolution kernels for volumetric data with the trick of vectorizing. Cesa et al. (2021) parameterizes filters with a band-limited basis to build E(3)-steerable kernels. Thomas et al. (2018) and Fuchs et al. (2020) extended the equivariance to graph neural networks.

In the context of robotics learning, Zhu et al. (2022) decouple rotation and translation symmetries to enable the robot to learn a planar grasp policy online within 1.5 h. Compared with Zhu et al. (2022) that formulated the planar grasping task as a bandit problem, our work focuses on pick-place tasks and learns from demonstrations. Wang et al. (2022) use SE(2) equivariance in Q learning to solve multi-step sequential manipulation pick-place tasks. Compared with Wang et al. (2022), our work leverages the larger  $SO(2) \times SO(2)$  symmetry group for the pick-conditioned place policy and tackles rearrangement tasks through the imitation learning (Hussein et al., 2017; Hester et al., 2018; Vecerik et al., 2017). Recently, various SE(3) equivariant architectures (Thomas et al., 2018; Fuchs et al., 2020; Chen et al., 2021; Deng et al., 2021) have been proposed and applied to solve manipulation problems. Simeonov et al. (2022) use Vector Neurons (Deng, Litany, Duan, Poulenc, Tagliasacchi, and Guibas, 2021) to get SE(3)-invariant object representations so that the model can manipulate objects in the same category with a few training demonstrations. Huang, Wang, Zhu, Walters and Platt (2022b) leverages the SE(3) invariance of the grasping evaluation function to enable better grasping performance. Xue et al. (2022) use SE(3)-equivariant key points to infer the object's pose for pick and place. However, most SE(3)-equivariant pick-place methods (Simeonov et al., 2022; Xue et al., 2022) require a segmentation model and a pre-trained point descriptor for each category, which limits their adaptations to various tasks. Although our proposed pick-place symmetry is defined on SE(2) in this work, we will briefly analyze how to extend the idea to SE(3)-pick-place problems in Proposition 3.

## 2. Background on symmetry groups

### 2.1. The groups $SO(2)$ and $C_n$

In this work, we primarily focus on rotations expressed by the group  $SO(2)$  and its cyclic subgroup  $C_n \subseteq SO(2)$ .  $SO(2)$  contains the continuous planar rotations  $\{\text{Rot}_\theta; 0 \leq \theta < 2\pi\}$ . The discrete subgroup  $C_n = \{\text{Rot}_\theta; \theta \in \{2\pi i/n | 0 \leq i < n\}\}$  contains only rotations by angles which are multiples of  $2\pi/n$ . The special Euclidean group  $SE(2) = SO(2) \times \mathbb{R}^2$  describes all translations and rotations of  $\mathbb{R}^2$ .

### 2.2. Representation of a group

A  $d$ -dimensional representation  $\rho: G \rightarrow GL_d$  of a group  $G$  assigns to each element  $g \in G$  an invertible  $d \times d$ -matrix  $\rho(g)$ . Different representations of  $SO(2)$  or  $C_n$  help to describe how different signals are transformed under rotations.

1. **The trivial representation**  $\rho_0: SO(2) \rightarrow GL_1$  assigns  $\rho_0(g) = 1$  for all  $g \in G$ , that is, no transformation under rotation.
2. **The standard representation**

$$\rho_1(\text{Rot}_\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

represents each group element by its standard rotation matrix. Notice that  $\rho_0$  and  $\rho_1$  can be used to represent elements from either  $\text{SO}(2)$  or  $C_n$ .

3. **The regular representation**  $\rho_{\text{reg}}$  of  $C_n$  acts on a vector in  $\mathbb{R}^n$  by cyclically permuting its coordinates  $\rho_{\text{reg}}(\text{Rot}_{2\pi/n})(x_0, x_1, \dots, x_{n-2}, x_{n-1}) = (x_{n-1}, x_0, x_1, \dots, x_{n-2})$ . We can rotate by multiples of  $2\pi/n$  by  $\rho_{\text{reg}}(\text{Rot}_{2\pi i/n}) = \rho_{\text{reg}}(\text{Rot}_{2\pi/n})^i$ .
4. **The quotient representation** of  $C_n$  for  $k$  dividing  $n$  is denoted  $\rho_{\text{quot}}^{C_n/C_k}$  and acts on  $\mathbb{R}^{n/k}$  by permuting  $|C_n|/|C_k|$  channels:  $\rho_{\text{quot}}^{C_n/C_k}(\text{Rot}_{2\pi i/n})(\mathbf{x})_j = (\mathbf{x})_{j+i \bmod (n/k)}$ , which implies features that are invariant under the action of  $C_k$ .
5. **The irreducible representation**  $\rho_{\text{irrep}}^i$  could be considered as the basis function with the order/frequency of  $i$ , such that any representation  $\rho$  of  $G$  could be decomposed as a *direct sum* of them:  $\rho(g) = Q^\top (\oplus_i \rho_{\text{irrep}}^i) Q$ , where  $Q$  is an orthogonal matrix.

For more details, we refer interesting readers to [Serre \(1977\)](#), [Weiler and Cesa \(2019\)](#), [Lang and Weiler \(2020\)](#), and [Cesa et al. \(2021\)](#).

### 2.3. Feature map transformations

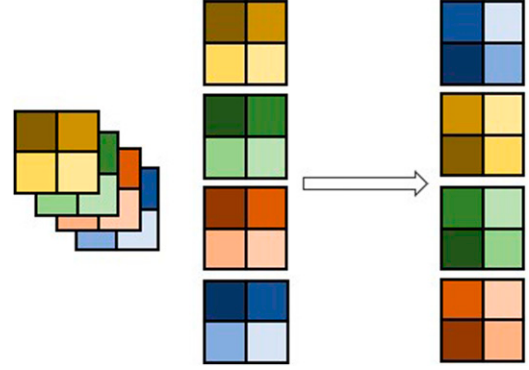
We formalize images and 2D feature maps as feature vector fields, that is, functions  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^c$ , which assign a feature vector  $f(\mathbf{x}) \in \mathbb{R}^c$  to each position  $\mathbf{x} \in \mathbb{R}^2$ . While in practice we discretize and truncate the domain of  $f$   $\{(i, j): 1 \leq i \leq W, 1 \leq j \leq W\}$ , here we will consider it to be continuous for the purpose of analysis. The action of an element  $g \in \text{SO}(2)$  on  $f$  is a combination of a rotation in the domain of  $f$  via  $\rho_1$  (this rotates the pixel positions) and a transformation in the channel space  $\mathbb{R}^c$  (aka. fiber space) by  $\rho \in \{\rho_0, \rho_1, \rho_{\text{reg}}, \rho_{\text{irrep}}\}$ . If  $\rho = \rho_{\text{reg}}$ , then the channels cyclically permute according to the rotation. If  $\rho = \rho_0$ , the channels do not change. We denote this action (the action of  $g$  on  $f$  via  $\rho$ ) by  $T_g^\rho(f)$ :

$$[T_g^\rho(f)](\mathbf{x}) = \rho(g) \cdot f(\rho_1(g)^{-1}\mathbf{x}). \quad (1)$$

For example, the action of  $T_g^{\rho_{\text{reg}}}(f)$  is illustrated in [Figure 2](#) for a rotation of  $g = \pi/2$  on a  $2 \times 2$  image  $f$  that uses  $\rho_{\text{reg}}$ . The expression  $\rho_1(g)^{-1}\mathbf{x}$  rotates the pixels via the standard representation. Multiplication by  $\rho(g) = \rho_{\text{reg}}(g)$  permutes the channels. For brevity, we will denote  $T_g^{\rho_{\text{reg}}} = T_g^{\rho_{\text{reg}}}$  and  $T_g^0 = T_g^{\rho_0}$ .

### 2.4. Equivariant mappings and steerable kernels

A function  $F$  is equivariant if it commutes with the action of the group,



**Figure 2.** Illustration of the action of  $T_g^{\text{reg}}$  on a  $2 \times 2$  image.

$$T_g^{\text{out}}(F(f)) = F\left(T_g^{\text{in}}(f)\right) \quad (2)$$

where  $T_g^{\text{in}}$  transforms the input to  $F$  by the group element  $g$  while  $T_g^{\text{out}}$  transforms the output of  $F$  by  $g$ . For example, if  $f$  is an image, then  $\text{SO}(2)$ -equivariance of  $F$  implies that it acts on  $f$  in the same way regardless of the orientation in which  $f$  is presented. That is, if  $F$  takes an image  $f$  rotated by  $g$  (RHS of equation (2)), then it is possible to recover the same output by evaluating  $F$  for the un-rotated image  $f$  and rotating its output (LHS of equation (2)). The most equivariant mappings between spaces of feature fields are *convolutions with  $G$ -steerable kernels* ([Weiler et al., 2018](#); [Jenner and Weiler, 2021](#)). Denote the input field type as  $\rho_{\text{in}}: G \rightarrow \mathbb{R}^{d_{\text{in}} \times d_{\text{in}}}$  and the output field type as  $\rho_{\text{out}}: G \rightarrow \mathbb{R}^{d_{\text{out}} \times d_{\text{out}}}$ . The  $G$ -steerable kernels are convolution kernels  $K: \mathbb{R}^n \rightarrow \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$  satisfying the *steerability constraint*, where  $n$  is the dimensionality of the space

$$K(g \cdot x) = \rho_{\text{out}}(g)K(x)\rho_{\text{in}}(g)^{-1} \quad (3)$$

## 3. Problem statement

This paper considers behavior cloning for planar pick and place problems. These problems are planar in the sense that the observation is a top-down image and the pick and place actions are motions to coordinates in the plane. Given a set of demonstrations that contains a sequence of one or more observation-action pairs  $(o_t, a_t)$ , the objective is to infer a policy  $p(a_t|o_t)$  where the action  $a_t = (a_{\text{pick}}, a_{\text{place}})$  describes both the pick and place components of action, and the observation  $o_t$  describes the current state in terms of a top-down image of the workspace.

Our model will encode this policy by factoring  $p(a_{\text{pick}}|o_t)$  and  $p(a_{\text{place}}|o_t, a_{\text{pick}})$  and representing them as two separate neural networks. This policy can be used to solve tasks that are solvable in a single time step (i.e., a single pick and place action) as well as tasks that require multiple pick and place actions to solve.  $a_{\text{pick}}$  and  $a_{\text{place}}$  are parameterized in terms of  $\text{SE}(2)$  coordinates  $(u, v, \theta)$ , where  $u, v$  denote the pixel coordinates of the gripper position and  $\theta$  denotes gripper

orientation.  $\theta_{\text{pick}}$  is defined with respect to the world frame and  $\theta_{\text{place}}$  is the delta action between the pick pose and place pose.

## 4. Transporter network

Before describing Equivariant Transporter Net, we analyze the original Transporter Net (Zeng et al., 2021) architecture from a different perspective.

### 4.1. Description of transporter net

Transporter Network (Zeng et al., 2021) solves the planar pick and place problem using the architecture shown in Figure 3. The pick network  $f_{\text{pick}}: o_t \mapsto p(u, v)$  maps and image  $o_t$  onto a probability distribution  $p(u, v)$  over pick position  $(u, v) \in \mathbb{R}^2$ . The output pick position  $a_{\text{pick}}^*$  is calculated by maximizing  $f_{\text{pick}}(o_t)$  over  $(u, v)$ . (Since Zeng et al. (2021) uses suction cups to pick, that work ignores pick orientation.) The place position and orientation is calculated as follows. First, an image patch  $c$  centered on  $a_{\text{pick}}^*$  is cropped from  $o_t$  to represent the pick action as well as the object. Then, the crop  $c$  is rotated  $n$  times to produce a stack of  $n$  rotated crops. We denote this stack of crops as

$$\mathcal{R}_n(c) = \left( T_{2\pi i/n}^0(c) \right)_{i=0}^{n-1}, \quad (4)$$

where we refer to  $\mathcal{R}_n$  as the “lifting” operator of  $C_n$ . Then,  $\mathcal{R}_n(c)$  is encoded using a neural network  $\psi$ . The original image,  $o_t$ , is encoded by a separate neural network  $\phi$ . The distribution over place location is evaluated by taking the cross-correlation between  $\psi$  and  $\phi$ ,

$$f_{\text{place}}(o_t, c) = \psi(\mathcal{R}_n(c)) \star \phi(o_t), \quad (5)$$

where  $\psi$  is applied independently to each of the rotated channels in  $\mathcal{R}_n(c)$ . Place position and orientation is calculated by maximizing  $f_{\text{place}}$  over the pixel position (for position) and the orientation channel (for orientation).

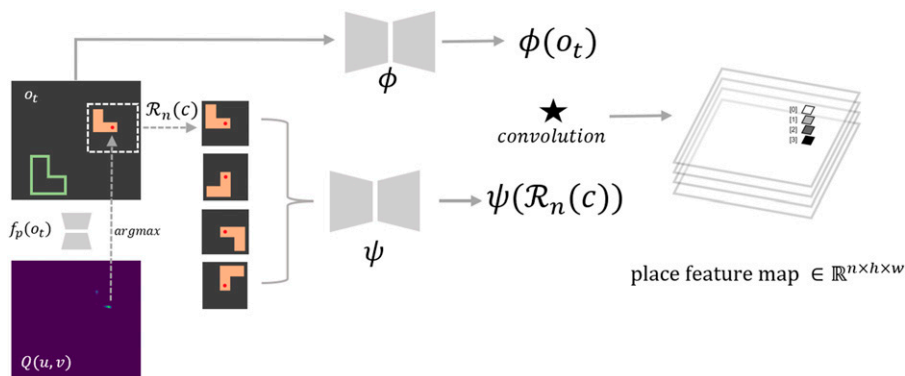


Figure 3. The architecture of transporter net.

### 4.2. Analysis of transporter net

The model architecture described above gives Transporter Network the following equivariance property.

**Proposition 1.** *The Transporter Net place network  $f_{\text{place}}$  is  $C_n$ -equivariant. That is, given  $g \in C_n$ , object image crop  $c$ , and scene image  $o_t$ ,*

$$f_{\text{place}}(o_t, T_g^0(c)) = \rho_{\text{reg}}(-g) f_{\text{place}}(o_t, c). \quad (6)$$

Proposition 1 expresses the following intuition. A rotation of  $g$  applied to the orientation of the object to be picked results in a  $-g$  change in the placing angle, which is represented by a permutation along the channel axis of the placing feature maps. We denote the permutation in the channel space as  $\rho_{\text{reg}}(-g)$ . This is a symmetry over the cyclic group  $C_n \subseteq \text{SO}(2)$  which is encoded directly into the model. It enables it to immediately generalize over different orientations of the object to be picked and thereby improves sample efficiency.

To prove Proposition 1, we start with some common lemmas. In order to understand continuous rotations of image data, it is helpful to consider a  $k$ -channel image as a mapping  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^k$  where the input  $\mathbb{R}^2$  defines the pixel space. We consider images centered at  $(0,0)$  and for non-integer values  $(x, y)$  we consider  $f(x, y)$  to be the interpolated pixel value. Similarly, let  $K: \mathbb{R}^2 \rightarrow \mathbb{R}^{l \times k}$  be a convolutional kernel where  $k$  is the number of the input channels and  $l$  is the number of the output channels. Although the input space is  $\mathbb{R}^2$ , we assume the kernel is  $r \times r$  pixels and  $K(x, y)$  is zero outside this set. The convolution can then be expressed by  $(K \star f)(\vec{v}) = \sum_{\vec{w} \in \mathbb{Z}^2} f(\vec{v} + \vec{w}) K(\vec{w})$ , where  $\vec{v} = (i, j) \in \mathbb{R}^2$ .

Without loss of generality, assume that  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  and define  $\tilde{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^n$  to be the  $n$ -fold duplication of  $f$  such that  $\tilde{f}(\vec{v}) = (f(\vec{v}), \dots, f(\vec{v}))$ . Consider a diagonal kernel  $\tilde{K}: \mathbb{R}^2 \rightarrow \mathbb{R}^{n \times n}$  where  $\tilde{K}(\vec{v})$  is a diagonal  $n \times n$  matrix  $\text{Diag}(K_1, \dots, K_n)$  and each  $K_i: \mathbb{R}^2 \rightarrow \mathbb{R}^{1 \times 1}$ .

For such inputs and kernels, we have the following permutation equivariance.

**Lemma 1.**

$$(\rho_{\text{reg}}(g)\tilde{K})\star\tilde{f} = \rho_{\text{reg}}(g)(\tilde{K}\star\tilde{f})$$

**Proof.** By definition  $h_i = (\tilde{K}\star\tilde{f})_i = K_i\star f$ . Define  $h = (h_1, \dots, h_n)$  and it is clear that permuting the  $1 \times 1$  kernels  $K_i$  also permutes  $h_i$ , so  $\rho_{\text{reg}}(g)h = (\rho_{\text{reg}}(g)\tilde{K})\star\tilde{f}$  as desired.

We require one more lemma on the equivariance of the lifting operator  $\mathcal{R}_n$ .

**Lemma 2.**

$$\mathcal{R}_n(T_g^0 f) = \rho_{\text{reg}}(-g)\mathcal{R}_n(f)$$

**Proof.** First, we compute

$$\mathcal{R}_n(f)(\vec{x}) = (f(\vec{x}), f(g^{-1}\vec{x}), \dots, f(g^{-(n-1)}\vec{x})).$$

Then both  $\mathcal{R}_n(T_g^0 f)$  and  $\rho_{\text{reg}}(-g)\mathcal{R}_n(f)$  equal to

$$(f(g^{-1}\vec{x}), \dots, f(g^{-(n-1)}\vec{x}), f(\vec{x})).$$

Proof of Proposition 1, we prove the  $C_n$ -place equivariance of Transporter Net under rotations of the picked object,

$$\psi(\mathcal{R}_n(T_g^0 c))\star\phi(o_i) = \rho_{\text{reg}}(-g)(\psi(\mathcal{R}_n(c))\star\phi(o_i)) \quad (7)$$

**Proof.** Since  $\psi$  is applied independently to each of the rotated channels in  $\mathcal{R}_n(c)$ , we denote  $\psi_n((f_1, \dots, f_n)) = (\psi(f_1), \dots, \psi(f_n))$ . By Lemma 2, the left-hand side of equation (7) is

$$\psi(\mathcal{R}_n(T_g^0 c))\star\phi(o_i) = \psi_n(\rho_{\text{reg}}(-g)\mathcal{R}_n(c))\star\phi(o_i).$$

Since  $\psi_n$  applies  $\psi$  on each component, it is equivariant to the permutation of components and thus the above equation becomes

$$= (\rho_{\text{reg}}(-g)\psi_n(\mathcal{R}_n(c))\star\phi(o_i).$$

Finally applying Lemma 1 gives

$$= \rho_{\text{reg}}(-g)(\psi_n(\mathcal{R}_n(c))\star\phi(o_i))$$

as desired.

The main idea of the proof is shown in Figure 4. Namely,  $\psi(\mathcal{R}_n(\cdot))$  is equivariant in the sense that rotating the crop  $c$  induces a cyclic shift in the channels of the output. Formally,  $\psi(\mathcal{R}_n(T_g^0 c)) = \rho_{\text{reg}}(-g)\psi(\mathcal{R}_n(c))$ . Noting that a permutation of the filters  $K$  in the convolution  $K\star\phi(o_i)$  induces the same permutation in the output feature maps completes the proof. Here,  $\psi$  is a simple CNN with no rotational equivariance. The equivariance results from the lifting operator  $\mathcal{R}_n$ .

However, only the place network of Transporter Net has the  $C_n$ -equivariance. Instead, our proposed method

incorporates not only the rotational equivariance in the pick network but also  $C_n \times C_n$ -equivariance in the place network.

## 5. Equivariant transporter

### 5.1. Equivariant pick

Our approach to the pick network is similar to that in Transporter Net Zeng et al. (2021) except that: (1) we explicitly encode *the pick symmetry* into the pick networks, thereby making pick learning more sample-efficient; (2) we consider the pick orientation so that we can use parallel-jaw grippers rather than just suction grippers.

*5.1.1. Model.* We propose an equivariant model for detecting the planar pick pose. First, we decompose the learning process of  $a_{\text{pick}} \in \text{SE}(2)$  into two parts,

$$p(a_{\text{pick}}) = p(u, v)p(\theta|u, v), \quad (8)$$

where  $p(u, v)$  denotes the probability of success when a pick exists at pixel coordinates  $u, v$  and  $p(\theta|u, v)$  is the probability that the pick at  $u, v$  should be executed with a gripper orientation of  $\theta$ . The distributions  $p(u, v)$  and  $p(\theta|u, v)$  are modeled as two neural networks:

$$f_p(o_i) \mapsto p(u, v), \quad (9)$$

$$f_\theta(o_i, (u, v)) \mapsto p(\theta|u, v). \quad (10)$$

Given this factorization, we can query the maximum of  $p(a_{\text{pick}})$  by evaluating  $(\hat{u}, \hat{v}) = \arg \max_{(u, v)} (p(u, v))$  and then  $\hat{\theta} = \arg \max_{\theta} (p(\theta|\hat{u}, \hat{v}))$ . This is illustrated in Figure 5. The bottom of Figure 5 shows the maximization of  $f_p^*$  at  $a_{\text{pick}}^*$ . The right side shows the evaluation of  $f_\theta$  for the image patch centered at  $a_{\text{pick}}^*$ .

*5.1.2. Pick symmetry.* There are two equivariant relationships that we would expect to be satisfied for planar picking:

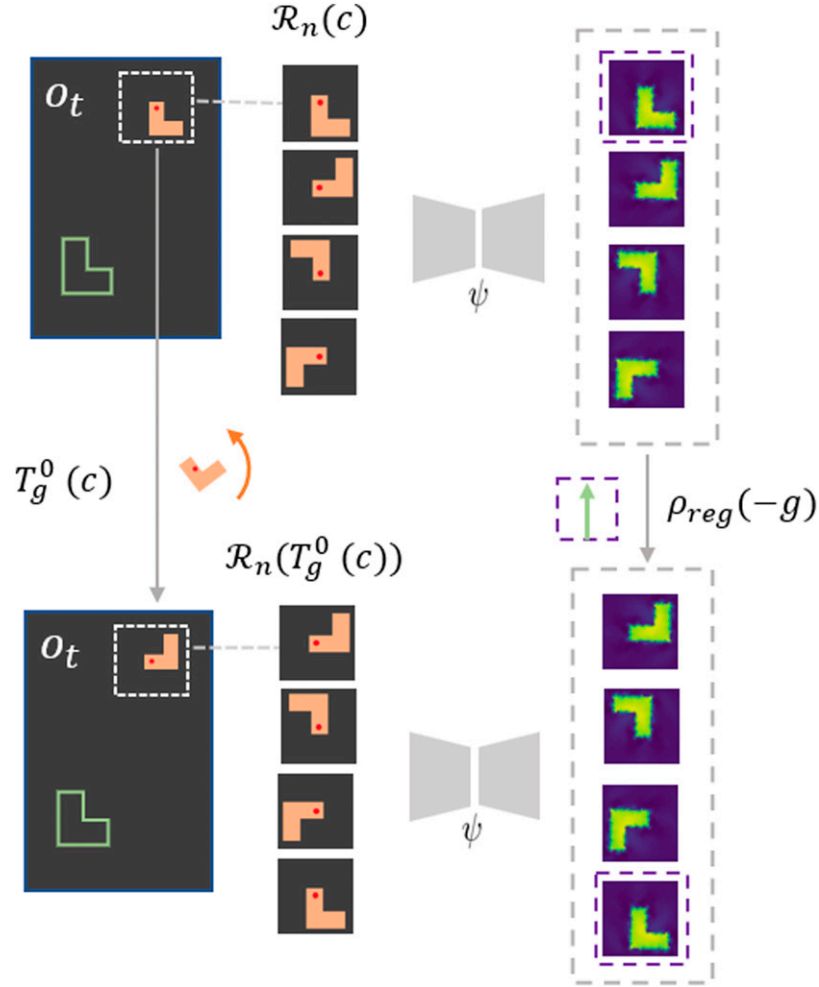
$$f_p(T_g^0(o_i)) = T_g^0(f_p(o_i)) \quad (11)$$

$$f_\theta(T_g^0(o_i), T_g^0(u, v)) = s(g)(f_\theta(o_i, (u, v))) \quad (12)$$

where  $s$  is the shift operator and satisfies  $s(g)f(x) = f(x + g)$ .

Equation (11) states that the pick location distribution found in an image rotated by  $g \in \text{SO}(2)$ , (LHS of equation (11)), should correspond to the distribution found in the original image subsequently rotated by  $g$ , (RHS of equation (11)).

Equation (12) says that the pick orientation distribution at the rotated grasp point  $T_g^0(u, v)$  in the rotated image  $T_g^0(o_i)$  (LHS of Equation (12)) should be shifted by  $g$  relative to the grasp orientation at the original grasp points in the original image (RHS of equation (12)).



**Figure 4.** Illustration of the main part of the proof of Proposition 1. Rotating the crop  $c$  induces a cyclic shift in the channels of the output  $\psi(\mathcal{R}_n(T_g^0 c)) = \rho_{\text{reg}}(-g)\psi(\mathcal{R}_n(c))$ .

We encode both  $f_p$  and  $f_\theta$  using equivariant convolutional layers (Weiler and Cesa, 2019) which constrain the models to represent only those functions that satisfy equations (11) and (12). Specifically, we select the trivial representation as the output type for  $f_p$  and the regular representation as the output type for  $f_\theta$ , which is a *special case*<sup>2</sup> of equation (12)

$$f_\theta(T_g^0(o_t), T_g^0(u, v)) = p_{\text{reg}}(g)(f_\theta(o_t, (u, v))) \forall g \in C_n \quad (13)$$

**5.1.3. Gripper orientation using the quotient group.** A key observation in planar picking is that, for many robots, the gripper is bilaterally symmetric, that is, grasp outcome is invariant when the gripper is rotated by  $\pi$ . We can encode this additional symmetry to reduce redundancy and save computational cost using the quotient group  $\text{SO}(2)/C_2$  which identifies orientations that are  $\pi$  apart. When using this quotient group for gripper orientation,  $s(g)$  in equation (12) is replaced with  $s(g \bmod \pi)$ <sup>3</sup> and  $\rho_{\text{reg}}$  in equation (13) is replaced with  $\rho_{\text{reg}}^{C_n/C_2}$ .

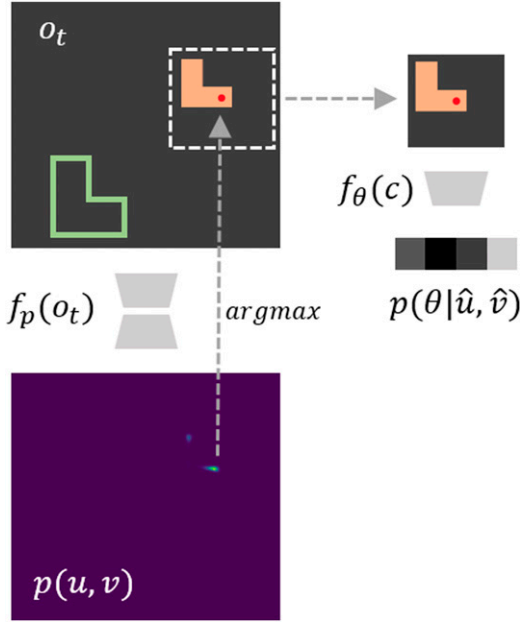
## 5.2. Equivariant place

Assumes that the object does not move during picking, given the picked object represented by the image patch  $c$  centered on  $a_{\text{pick}}$ , the place network models the distribution of  $a_{\text{place}} = (u_{\text{place}}, v_{\text{place}}, \theta_{\text{place}})$  by

$$f_{\text{place}}(o_t, c) \mapsto p(a_{\text{place}} | o_t, a_{\text{pick}}), \quad (14)$$

where  $p(a_{\text{place}} | o_t, a_{\text{pick}})$  denotes the probability that the object at  $a_{\text{pick}}$  in scene  $o_t$  should be placed at  $a_{\text{place}}$ .

Our place model architecture closely follows that of Transporter Net (Zeng et al., 2021). The main difference is that we explicitly encode equivariance constraints on both  $\phi$  and  $\psi$  networks. As a result of this change: (1) we are able to simplify the model by transposing the lifting operation  $\mathcal{R}_n$  and the processing by  $\phi$ ; (2) our new model is equivariant with respect to a larger symmetry group  $C_n \times C_n$ , compared to Transporter Net which is only equivariant over  $C_n$ .



**Figure 5.** Equivariant Transporter Pick model. First, we find the pick position  $a_{pick}^*$  by evaluating the  $\text{argmax}$  over  $f_p(o_t)$ . Then, we evaluate  $f_\theta$  for the image patch centered on  $a_{pick}^*$ .

**5.2.1. Equivariant  $\phi$  and  $\psi$ .** We explicitly encode both  $\phi$  and  $\psi$  as equivariant models that satisfy the following constraints:

$$\phi\left(T_g^0(o_t)\right) = T_g^0(\phi(o_t)) \quad (15)$$

$$\psi\left(T_g^0(c)\right) = T_g^0(\psi(c)) \quad (16)$$

for  $g \in \text{SO}(2)$ . The equivariance constraint of equation (15) says that when the input image rotates, we would expect the place location to rotate correspondingly. This constraint helps the model generalize across place orientations. The constraint of equation (16) says that when the picked object rotates (represented by the image patch  $c$ ), then the place orientation should correspondingly rotate.

**5.2.2. Place model.** When the equivariance constraint of equation (16) is satisfied, we can exchange  $\mathcal{R}_n$  (the lifting operation) with  $\psi$ :

$$\psi(\mathcal{R}_n(c)) = \mathcal{R}_n(\psi(c))$$

This equality is useful because it means that we only need to evaluate  $\psi$  for one image patch and rotate the feature map rather than processing the stack of image patches  $\mathcal{R}_n(c)$ —something that is computationally cheaper. The resulting place model is then:

$$f'_{\text{place}}(o_t, c) = \mathcal{R}_n(\psi(c)) \star \phi(o_t) \quad (17)$$

$$= \Psi(c) \star \phi(o_t), \quad (18)$$

where Equation section 18 substitutes  $\Psi(c) = \mathcal{R}_n(\psi(c))$  to simplify the expression. Here, we use  $f'_{\text{place}}$  to denote Equivariant Transporter Net defined using equivariant  $\phi$  and  $\psi$  in contrast to the baseline Transporter Net  $f_{\text{place}}$  of equation (5). Note that both  $f_{\text{place}}$  and  $f'_{\text{place}}$  satisfy Proposition 1. However,  $f_{\text{place}}$  accomplishes this by symmetrizing a non-equivariant network (i.e., evaluating  $\psi(\mathcal{R}_n(c))$ ) whereas our model  $f'_{\text{place}}$  encodes the symmetry directly into  $\psi$ .

### 5.3. Place symmetry of the equivariant transporter network

**5.3.1.  $C_n \times C_n$ -place symmetry.** As Proposition 1 demonstrates, the baseline Transporter Net model (Zeng et al., 2021) encodes the symmetry that rotations of the object to be picked (represented by  $c$ ) should result in corresponding rotations of the place orientation for that object. However, the pick-conditioned place has a second symmetry that is not encoded in Transporter Net: rotations of the placement (represented by  $o_t$ ) should also result in corresponding rotations of the place orientation. In fact, as we demonstrate in Proposition 2 below, we encode this *second type* of symmetry by enforcing the constraints of equations (15) and (16). Essentially, we go from the  $C_n$ -place symmetric model to a  $C_n \times C_n$ -place symmetric model.

**Proposition 2.** *Equivariant Transporter Net  $f'_{\text{place}}$  is  $C_n \times C_n$ -equivariant. That is, given rotations  $g_1 \in C_n$  of the picked object and  $g_2 \in C_n$  of the scene, we have that*

$$f'_{\text{place}}\left(T_{g_1}^0(c), T_{g_2}^0(o_t)\right) = \rho_{\text{reg}}(g_2 - g_1) T_{g_2}^0 f'_{\text{place}}(c, o_t). \quad (19)$$

Proposition 2 is illustrated in Figure 6. The top of Figure 6 going left to right shows the rotation of both the object by  $g_1$  (in orange) and the place pose by  $g_2$  (in green). The LHS of equation (19) evaluates  $f'_{\text{place}}$  for these two rotated images. The lower left of Figure 6 shows  $f'_{\text{place}}(c, o_t)$ . Going left to right at the bottom of Figure 6 shows the pixel-rotation by  $T_{g_2}^0$  and the channel permutation by  $g_2 - g_1$  (RHS of equation (19)).

To prove Proposition 2, we introduce one more lemma.

**Lemma 3.**

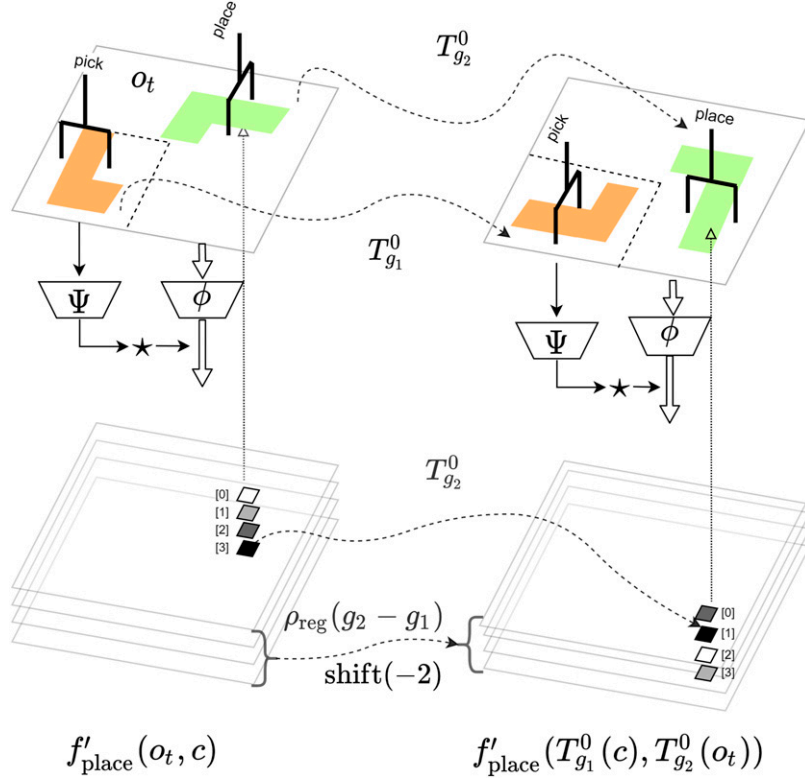
$$\left(T_g^0(K \star f)\right)(\vec{v}) = \left(\left(T_g^0 K\right) \star \left(T_g^0 f\right)\right)(\vec{v}) \quad (20)$$

**Proof.** We evaluate the left-hand side of equation (20):

$$T_g^0(K \star f)(\vec{v}) = \sum_{\vec{w} \in \mathbb{Z}^2} f(g^{-1}\vec{v} + \vec{w})K(\vec{w}).$$

Re-indexing the sum with  $\vec{y} = g\vec{w}$ ,

$$= \sum_{\vec{y} \in \mathbb{Z}^2} f(g^{-1}\vec{v} + g^{-1}\vec{y})K(g^{-1}\vec{y})$$



**Figure 6.** Equivariance of our placing network under the rotation of the object and the placement. A  $\pi/2$  rotation on  $c$  and a  $-\pi/2$  rotation on  $o_t$  are equivariant to: i), a  $-\pi/2$  rotation on the placing location, and ii), the shift on the channel of placing rotation angle from  $3\pi/2$  (the last channel) to  $\pi/2$  (the second channel).

is by definition

$$\begin{aligned} &= \sum_{\vec{y} \in \mathbb{Z}^2} (T_g^0 f)(\vec{v} + \vec{y}) (T_g^0 K)(\vec{y}) \\ &= ((T_g^0 K) \star (T_g^0 f))(\vec{v}) \end{aligned}$$

as desired.

**Proof of Proposition 2** Recall  $\Psi(c) = \psi(\mathcal{R}_n(c))$ . We now prove Proposition 2,

$$\begin{aligned} &\Psi(T_{g_1}^0(c)) \star \phi(T_{g_2}^0(o_t)) = \\ &\rho_{\text{reg}}(g_2 - g_1) (T_{g_2}^0(\Psi(c) \star \phi(o_t))) \end{aligned}$$

**Proof.** We first prove the equivariance under rotations of the placement  $o_t$ . We claim

$$\Psi(c) \star \phi(T_g^0(o_t)) = T_g^{\text{reg}}(\Psi(c) \star \phi(o_t)). \quad (21)$$

Evaluating the left-hand side of equation (21),

$$\begin{aligned} &\Psi(c) \star \phi(T_g^0(o_t)) \\ &= \Psi(c) \star T_g^0 \phi(o_t) \text{ (equivariance of } \phi) \\ &= (T_g^0 T_{g^{-1}}^0 \Psi(c)) \star (T_g^0 \phi(o_t)) \\ &= T_g^0 (T_{g^{-1}}^0 \Psi(c) \star \phi(o_t)) \text{ (Lemma 3)} \\ &= T_g^0 (T_{g^{-1}}^0 \mathcal{R}_n(\psi(c)) \star \phi(o_t)) \\ &= T_g^0 (\mathcal{R}_n(T_{g^{-1}}^0 \psi(c)) \star \phi(o_t)) \text{ (equivariance of } \mathcal{R}_n) \\ &= T_g^0 (\mathcal{R}_n(\psi(T_{g^{-1}}^0 c)) \star \phi(o_t)) \text{ (equivariance of } \psi) \\ &= T_g^0 ((\rho_{\text{reg}}(g) \Psi(c)) \star \phi(o_t)) \text{ (Lemma 2)} \\ &= T_g^0 \rho_{\text{reg}}(g) (\Psi(c) \star \phi(o_t)) \text{ (Lemma 1)} \\ &= T_g^{\text{reg}}(\Psi(c) \star \phi(o_t)). \end{aligned}$$

In the last step,  $T_g^{\text{reg}} = \rho_{\text{reg}}(g) T_g^0 = T_g^0 \rho_{\text{reg}}(g)$  since  $T_g^0$  and  $\rho_{\text{reg}}(g)$  commute as  $\rho_{\text{reg}}(g)$  acts on the channel space and  $T_g^0$  acts on the base space. This proves the claim of equation (21).

Recall  $\Psi(c) = \mathcal{R}_n(\psi(c))$ . Using the equivariance of  $\psi$ , Proposition 1 could be reformulated as

$$\Psi\left(T_g^0 c\right) \star \phi(o_t) = \rho_{\text{reg}}(-g)(\Psi(c) \star \phi(o_t)) \quad (22)$$

Evaluating the left-hand side of equation (22),

$$\begin{aligned} & \Psi\left(T_g^0 c\right) \star \phi(o_t) \\ &= \mathcal{R}_n\left(\psi\left(T_g^0 c\right)\right) \star \phi(o_t) \quad (\Psi(c) = \mathcal{R}_n(\psi(c)) \text{ by def.}) \\ &= \psi\left(\mathcal{R}_n\left(T_g^0 c\right)\right) \star \phi(o_t) \quad (\text{equivariance of } \psi) \\ &= \rho_{\text{reg}}(-g)(\psi(\mathcal{R}_n(c)) \star \phi(o_t)) \quad (\text{Proposition 1}) \\ &= \rho_{\text{reg}}(-g)(\mathcal{R}_n(\psi(c)) \star \phi(o_t)) \quad (\text{equivariance of } \psi) \\ &= \rho_{\text{reg}}(-g)(\Psi(c) \star \phi(o_t)) \end{aligned}$$

Combining equation (21) with equation (22) realizes the Proposition 2.

**5.3.2. Translational symmetry.** Note that in addition to the two rotational symmetries enforced by our model, it also has translational symmetry. Since the rotational symmetry is realized by additional restrictions to the weights of kernels of convolutional networks, the rotational symmetry is in addition to the underlying shift equivariance of the convolutional network. Thus, the full symmetry group enforced is the group generated by  $C_n \times C_n \times (\mathbb{R}^2, +)$ . Equivariant neural networks learn effectively on a lower dimensional space, the equivalence classes of samples under the group action.

**5.3.3. From  $C_n \times C_n$ -place symmetry to  $SO(2) \times SO(2)$ .** The above place symmetry is limited to the cyclic group due to the role of  $\mathcal{R}_n$ , though as  $n \rightarrow \infty$ ,  $C_n$  equals  $SO(2)$ . We show the generalization of the  $C_n \times C_n$ -place symmetry and  $SO(2) \times SO(2)$  place symmetry below.

Given  $g \in G$  for  $G \subseteq SO(2)$ , an equivariant model  $\phi$  satisfying  $\phi(T_g^0(o_t)) = T_g^0(\phi(o_t))$  and a function  $\bar{\Psi}: c \rightarrow K$  satisfying the equivariant constraint  $\bar{\Psi}(T_g^0 c) = T_g^0 \bar{\Psi}(c)$ , where  $c$  is the crop  $\in \mathbb{R}^2$  and  $K: \mathbb{R}^2 \rightarrow \mathbb{R}^{d_{\text{out}} \times d_{\text{trivial}}}$  is a 2D steerable kernel with trivial representation as the input type. The cross-correlation between  $\bar{\Psi}(c)$  and  $\phi(T_g^0(o_t))$  satisfies

$$\bar{\Psi}\left(T_{g_1}^0(c)\right) \star \phi\left(T_{g_2}^0(o_t)\right) = \rho_{\text{out}}(g_2 - g_1)\left(T_{g_2}^0(\bar{\Psi}(c) \star \phi(o_t))\right) \quad (23)$$

Proposition 3 states that to satisfy the cross-type place symmetry, one necessary condition is that the output of  $\bar{\Psi}$  is a steerable kernel. It generalizes Proposition 2 to either  $C_n$  or  $SO(2)$ . In fact,  $\Psi(c) = \mathcal{R}_n(\psi(c))$  combining the lift operator  $\mathcal{R}_n$  and the equivariant constraint of  $\psi$  shown in equation (16) is a special case of  $\bar{\Psi}(c)$ .  $\mathcal{R}_n: \mathbb{R}^2 \rightarrow K$  outputting a steerable kernel  $K$  that takes the regular representation of  $C_n$  as the output type and satisfies the

steerability constraint of equation (3). When using *irreducible representations* as the output type, we can instantiate  $\rho_{\text{out}}(g_2 - g_1)$  in RHS of equation (23) as  $\rho_{\text{irrep}}(g_2 - g_1)$  which is equivalent to  $s(g_2 - g_1)$  after Inverse Fourier Transform.

To prove Proposition 3, we first introduce another lemma.

**Lemma 4.** A 2D steerable kernel  $K: \mathbb{R}^2 \rightarrow \mathbb{R}^{d_{\text{out}} \times d_{\text{trivial}}}$  satisfies

$$T_g^0 K(x) = \rho_{\text{out}}(g^{-1})K(x) \quad (24)$$

**Proof.** Recall that  $\rho_0(g)$  is an identity mapping. Substituting  $\rho_{\text{in}}$  with  $\rho_0(g)$  and  $g^{-1}$  with  $g$  in the steerability constraint  $K(g \cdot x) = \rho_{\text{out}}(g)K(x)\rho_{\text{in}}(g)^{-1}$  shown in equation (3) completes the proof.

$$\begin{aligned} T_g^0 K(x) &= K(g^{-1}x) \\ &= \rho_{\text{out}}(g^{-1})K(x)\rho_{\text{in}}(g) \\ &= \rho_{\text{out}}(g^{-1})K(x) \end{aligned}$$

Lemma 4 states that when the input type is the trivial representation, a spatial rotation of the steerable kernel is the same as the inverse channel space transformation. With Lemma 4 in hand, we start the proof of Proposition 3

**Proof.** Similar to the proof of Proposition 2, we first show the equivariance under rotations of the placement  $o_t$ . We claim

$$\bar{\Psi}(c) \star \phi\left(T_g^0(o_t)\right) = T_g^{\text{out}}(\bar{\Psi}(c) \star \phi(o_t)) \quad (25)$$

Starting from the left-hand side of equation (25),

$$\begin{aligned} & \bar{\Psi}(c) \star \phi\left(T_g^0(o_t)\right) \\ &= \bar{\Psi}(c) \star T_g^0 \phi(o_t) \quad (\text{equivariance of } \phi) \\ &= \left(T_g^0 T_{g^{-1}}^0 \bar{\Psi}(c)\right) \star \left(T_g^0 \phi(o_t)\right) \\ &= T_g^0 \left(T_{g^{-1}}^0 \bar{\Psi}(c) \star \phi(o_t)\right) \quad (\text{Lemma 3}) \\ &= T_g^0 (\rho_{\text{out}}(g) \bar{\Psi}(c) \star \phi(o_t)) \quad (\text{Lemma 4}) \\ &= T_g^{\text{out}}(\bar{\Psi}(c) \star \phi(o_t)). \end{aligned}$$

Then, we propose the equivariance under rotations of the picked object as

$$\bar{\Psi}\left(T_g^0(c)\right) \star \phi(o_t) = \rho_{\text{out}}(-g)(\bar{\Psi}(c) \star \phi(o_t)) \quad (26)$$

Evaluating the left-hand side of equation (26),

$$\begin{aligned} & \bar{\Psi}\left(T_g^0(c)\right) \star \phi(o_t) \\ &= T_g^0 \bar{\Psi}(c) \star \phi(o_t) \quad (\text{equivariance of } \bar{\Psi}) \\ &= \rho_{\text{out}}(-g) \bar{\Psi}(c) \star \phi(o_t) \quad (\text{Lemma 4}) \end{aligned}$$

Combining equation (25) with equation (26) realizes the Proposition 3.

Note that Proposition 3 gives the way to realize the SO(2) version of our model and provide some insights to extend it to 3D signals without limitations. That is, generating the 3D dynamic steerable kernels from the crop signal. But in this work, we primarily focus on the discrete  $C_n$  group since it is easy to compare with the baseline *Transporter Net* on Ravens-10 Benchmark.

#### 5.4. Analyzing equivariance under Proposition 2

We summarize some important properties from the larger symmetry group of our place network and provide an intuitive explanation for each one. Recall that Proposition 2 states

$$\begin{aligned} \Psi\left(T_{g_1}^0(c)\right) \star \phi\left(T_{g_2}^0(o_t)\right) \\ = \rho_{\text{reg}}(g_2 - g_1) \left( T_{g_2}^0(\Psi(c) \star \phi(o_t)) \right). \end{aligned}$$

Then we have the following properties:

**5.4.1. Equivariance property.** Setting either  $g_1 = 0$  or  $g_2 = 0$  we get, respectively,

$$\Psi\left(T_g^0(c)\right) \star \phi(o_t) = \rho_{\text{reg}}(-g) (\Psi(c) \star \phi(o_t)) \quad (27)$$

$$\Psi(c) \star \phi\left(T_g^0(o_t)\right) = T_g^{\text{reg}}(\Psi(c) \star \phi(o_t)) \quad (28)$$

These show the equivariance of our network  $f_{\text{place}}$  under either a rotation  $g \in C_n$  of the object or the placement.

**5.4.2. Invariance property.** Setting  $g_1 = g_2$ , we get

$$\Psi\left(T_g^0(c)\right) \star \phi\left(T_g^0(o_t)\right) = T_g^0(\Psi(c) \star \phi(o_t)). \quad (29)$$

This equation demonstrates that a rotation  $g$  on the whole observation  $o_t$  including the objects does not change the placing angle but rotates the placing location by  $g$ . Although data augmentation could help non-equivariant models learn this property, our networks observe it by construction. Note that for the discrete group, data augmentation propagates to every element within the group.

**5.4.3. Relativity property.** Related to equation (27), we also have

$$\Psi\left(T_g^0(c)\right) \star \phi(o_t) = \rho_{\text{reg}}(-g) \left( T_g^0(\Psi(c) \star \phi(T_{-g}^0(o_t))) \right)$$

This equation defines the *dual* relationship between a rotation on  $c$  by  $g$  and an inverse rotation  $-g$  on  $o_t$ . Intuitively,  $c$  could be considered as the L-shaped block and  $o_t$  can be regarded as the L-shaped slot. A rotation on the picked object is equivariant to an inverse rotation on the placement under some transformation.

#### 5.5. Goal-conditioned equivariant transporter network

The goal-conditioned pick-place task is an important branch in learning manipulation skills where the goal could be represented as language instructions, images, or other customized definitions. Seita et al. (2021) extended Transporter Net to solve image-based goal-conditioned tasks. In this setting, the goal is represented explicitly as an image that is part of the problem input rather than implicitly as part of the observation. Two goal-conditioned architectures are proposed in Seita et al. (2021). *Transporter-Goal-Stack* stacks the current  $o_t$  and goal  $o_g$  images channel-wise and passes it as input through a standard Transporter Network. *Transporter-Goal-Split* processes the goal image  $o_g$  through a separate Fully Convolution Network  $\phi_{\text{goal}}$  to generate dense features to be combined with dense features of  $\phi_{\text{query}}(o_t)$  using the Hadamard product to infer the goal-conditioned pick

$$f_{\text{query}} = \phi_{\text{query}}(o_t) \odot \phi_{\text{goal}}(o_g) \quad (30)$$

and evaluate the goal-conditioned place with

$$f_{\text{key}} = \phi_{\text{key}}(o_t) \odot \phi_{\text{goal}}(o_g) \quad (31)$$

$$p(a_{\text{place}} | o_t, o_g, a_{\text{pick}}) = \mathcal{R}_n(f_{\text{query}}[a_{\text{pick}}]) \star f_{\text{key}} \quad (32)$$

where  $f_{\text{query}}[a_{\text{pick}}]$  denotes the crop of the dense feature map  $f_{\text{query}}$  centered on  $a_{\text{pick}}$  and  $\odot$  is the Hadamard product.

Since the pick-place symmetries also exist in goal-conditioned tasks, we realize the goal-conditioned equivariant transporter with some simple modifications. Denote  $\parallel$  as the channel-wise concatenation, the  $C_n$ -equivariance of the picking network holds when stacking  $o_t$  and  $o_g$  as the input:

$$f_p\left(T_g^0(o_t \parallel o_g)\right) = T_g^0 f_p(o_t \parallel o_g) \quad (33)$$

The  $C_n \times C_n$ -equivariance of the placing model also holds

$$\begin{aligned} \Psi\left(T_{g_1}^0((o_t \parallel o_g)[a_{\text{pick}}])\right) \star \phi\left(T_{g_2}^0(o_t \parallel o_g)\right) = \\ \rho_{\text{reg}}(g_2 - g_1) \left( T_{g_2}^0(\Psi((o_t \parallel o_g)(a_{\text{pick}})) \star \phi(o_t \parallel o_g)) \right) \end{aligned} \quad (34)$$

Based on the two equations above, we define *Equivariant-Transporter-Goal-Stack* to solve goal-conditioned tasks.

#### 5.6. Model architecture details

**5.6.1. Pick model  $f_p$  (equation (9)).** The input to  $f_p$  is a 4-channel RGB-D image  $o_t \in \mathbb{R}^{4 \times H \times W}$ . The output is a feature map  $p(u, v) \in \mathbb{R}^{H \times W}$  which encodes a distribution over pick location.  $f_p$  is implemented as an 18-layer equivariant residual network with a U-Net (Ronneberger et al., 2015) as the main block. The U-net has eight residual blocks (each block contains two equivariant convolution layers (Weiler

and Cesa, 2019) and one skip connection): four residual blocks (He et al., 2016) are used for the encoder and the other four residual blocks are used for the decoder. The encoding process trades spatial dimensions for channels with max-pooling in each block; the decoding process upsamples the feature embedding with bilinear-upsampling operations. The first layer maps the trivial representation of  $o_t$  to regular representation and the last equivariant layer transforms the regular representation back to the trivial representation, followed by image-wide softmax. ReLU activations (Nair and Hinton, 2010) are interleaved inside the network.

**5.6.2. Pick model  $f_\theta$  (equation (10)).** Given the picking location  $(u^*, v^*)$ , the pick angle network  $f_\theta$  takes as input a crop  $c \in \mathbb{R}^{4 \times H_1 \times W_1}$  centered on  $(u^*, v^*)$  and outputs the distribution  $p(\theta|u, v) \in \mathbb{R}^{n/2}$ , where  $n$  is the size of the rotation group (i.e.,  $n = |C_n|$ ). The first layer maps the trivial representation of  $c$  to a quotient regular representation followed by three residual blocks containing max-pooling operators. This goes to two equivariant convolution layers and then to an average pooling layer.

**5.6.3. Place models  $\phi$  and  $\psi$ .** Our place model has two equivariant convolution networks,  $\phi$  and  $\psi$ , and both have similar architectures to  $f_p$ . The network  $\phi$  takes as input a zero-padded version of the 4-channel RGBD observation  $o_t$ ,  $\text{pad}(o_t) \in \mathbb{R}^{4 \times (H+d) \times (W+d)}$ , and generates a dense feature map,  $\phi(\text{pad}(o_t)) \in \mathbb{R}^{(H+d) \times (W+d)}$ , where  $d$  is the padding size. The network  $\psi$  takes as input the image patch  $c \in \mathbb{R}^{4 \times H_2 \times W_2}$  and outputs  $\psi(c) \in \mathbb{R}^{H_2 \times W_2}$ . After applying rotations of  $C_n$  to  $\psi(c)$ , the transformed dense embeddings  $\Psi(c) \in \mathbb{R}^{n \times H_2 \times W_2}$  are cross-correlated with  $\phi(\text{pad}(o_t))$  to generate the placing action distribution  $p(a_{\text{place}}|o_t, a_{\text{pick}}) \in \mathbb{R}^{n \times H \times W}$ , where the channel axis  $n$  corresponds to placing angles,  $2\pi i/n$  for  $0 \leq i < n$ .

**5.6.4. Group types and sizes.** The networks  $f_p$ ,  $\psi$ , and  $\phi$ :  $\rho_0 \rightarrow \rho_0$ , which are all defined using  $C_6$  regular representations in the intermediate layers. The ablation study of the group size of the latent feature is discussed in the Experiment section. The network  $f_\theta$ :  $\rho_0 \rightarrow \rho_{\text{quot}}$  is defined using the quotient representation  $C_{36}/C_2$ , which corresponds to the number of allowed pick orientations. The lift operator  $\mathcal{R}_n$  is implemented with  $C_{36}$  cyclic group, which allows 36 different place orientations. Both the number of allowed pick and place orientations are hyper-parameters and could be selected flexibly based on the task precision. Our choice of the  $\pi/18$  discretization, that is, 18 bilateral-symmetric pick orientations and 36 place orientations, follows the settings of Ravens-10 benchmark (Zeng et al., 2021).

**5.6.5. Training details.** We train Equivariant Transporter Network with the Adam (Kingma and Ba, 2014) optimizer with a fixed learning rate of  $10^{-4}$ . It takes about  $0.8 \text{ s}^4$  to complete one SGD step with a batch size of one on an

NVIDIA Tesla V100 SXM2 GPU. Compared with the baseline transporter net which takes around 0.6 s to complete one SGD step on the same setting, the equivariant constraint on the weight updating increases 33% computation load. In fact, Equivariant Transporter Net converges faster than the baseline Transporter Net as shown in Figure 9. This is due to that the larger symmetry group results in a smaller dimensional sample space and thus better coverage by the training data. For each task, we train a single-policy network and evaluate the performance every 1k steps on 100 unseen tests. On most tasks, the best performance is achieved in less than 10k SGD steps.

## 6. Experiments

We evaluate Equivariant Transporter using the Ravens-10 Benchmark (Zeng et al., 2021) and our variations thereof.

### 6.1. Tasks

**6.1.1. Ravens-10 tasks.** Ravens-10 is a behavior cloning simulation benchmark for manipulation, where each task owns an oracle that can sample expert demonstrations from the distribution of successful picking and placing actions with access to the ground-truth pose of each object. The 10 tasks of Ravens can be classified into three categories: *Single-object manipulation tasks* (block-insertion, align-box-corner); *multiple-object manipulation tasks* (place-red-in-green, towers-of-Hanoi, stack-block-pyramid, palletizing-boxes, assembling-kits, packing-boxes); *deformable-object manipulation task* (manipulating-rope, sweeping-piles).

Here we provide a short description of Ravens-10 Environment, we refer readers to Zeng et al. (2021) for details. The poses of objects and placements in each task are randomly sampled in the workspace without collision. Performance on each task is evaluated in one of two ways: (1) *pose*: translation and rotation error relative to target pose is less than a threshold  $\tau = 1 \text{ cm}$  and  $\omega = \pi/12$ , respectively. Tasks: block-insertion, towers-of-Hanoi, place-red-in-green, align-box-corner, stack-block-pyramid, assembling-kits. Partial scores are assigned to multiple-action tasks. (2) *Zone*: Ravens-10 discretizes the 3D bounding box of each object into  $2 \text{ cm}^3$  voxels. The Total reward is calculated by # of voxels in target zone/total # of voxels. Tasks: palletizing-boxes, packing-boxes, manipulating-cables, sweeping-piles. Note that pushing objects could also be parameterized with  $a_{\text{pick}}$  and  $a_{\text{place}}$  that correspond to the starting pose and the ending pose of the end effector.

1. **Block-insertion:** Picking up an L-shape block and placing it into an L-shaped fixture.
2. **Place-red-in-green:** picking up red cubes and placing them into green bowls. There could be multiple bowls and cubes with different colors.

3. **Towers-of-Hanoi:** sequentially picking up disks and placing them into pegs such that all three disks initialized on one peg are moved to another, and that only smaller disks can be on top of larger ones.
4. **Align-box-corner:** picking up a randomly sized box and placing it to align one of its corners to a green L-shaped marker labeled on the tabletop. This task requires precision and generalization ability to new box sizes.
5. **Stack-block-pyramid:** sequentially picking up six blocks and stacking them into a pyramid of 3-2-1.
6. **Palletizing-boxes:** picking up 18 boxes and stacking them on top of a pallet.
7. **Assembling-kits:** picking five shaped objects (randomly sampled with replacement from a set of 20 objects, where 14 objects are used during training and six objects are held out for testing) and fitting them to corresponding silhouettes of the objects on a board. This task requires generalizing to new objects.
8. **Packing-boxes:** picking and placing randomly sized boxes tightly into a randomly sized container.
9. **Manipulating-rope:** manipulating a deformable rope such that it connects the two endpoints of an incomplete 3-sided square (colored in green).
10. **Sweeping-piles:** pushing piles of small objects (randomly initialized) into a desired target goal zone on the tabletop marked with green boundaries. The task is implemented with a pad-shaped end effector.

## 6.2. Ravens-10 tasks modified for the parallel-jaw gripper

We select five tasks (block-insertion, align-box-corner, place-red-in-green, stack-block-pyramid, palletizing-boxes) from Ravens-10 and replaced the suction cup with the Franka Emika gripper, which requires additional picking angle inference. [Figure 7](#) illustrates the initial state and completion state for each of these five tasks. For each of these five tasks, we defined an oracle agent. Since the Transporter Net framework assumes that the object does not move during picking, we defined these expert generators such that this was the case.

**6.2.1. Goal-conditioned tasks.** We design four image-based goal-conditioned tasks (goal-conditioned block-insertion, goal-conditioned block-pyramid, goal-conditioned four blocks-no, goal-conditioned cable-align) based on ravens-10, as shown in [Figure 8](#).

For each of the four tasks, objects are generated with random poses on the workspace and there is no placement in the observation  $o_t$ . The robot must use pick-place actions to manipulate the objects to the target pose specified in the goal images. For the goal-conditioned cable-align task, the robot needs to align the rope to the straight line shown in the goal image.

## 6.3. Training and evaluation

**6.3.1. Training.** For each task, we produce a dataset of  $n$  expert demonstrations, where each demonstration contains a sequence of one or more observation-action pairs  $(o_t, \bar{a}_t)$  (or  $(o_t, o_g, \bar{a}_t)$  for goal-conditioned tasks). Each action  $\bar{a}_t$  contains an expert picking action  $\bar{a}_{\text{pick}}$  and an expert placing action  $\bar{a}_{\text{place}}$ . We use  $\bar{a}_t$  to generate one-hot pixel maps as the ground-truth labels for our picking model and placing model. The models are trained using a cross-entropy loss.

**6.3.2. Metrics.** We measure performance the same way as it was measured in Transporter Net ([Zeng et al., 2021](#))—using a metric in the range of 0 (failure) to 100 (success). Partial scores are assigned to multiple-action tasks. For example, in the block-stacking task where the agent needs to construct a 6-block-pyramid, each successful rearrangement is credited with a score of 16.667. We report the highest validation performance during training, averaging over 100 unseen tests for each task.

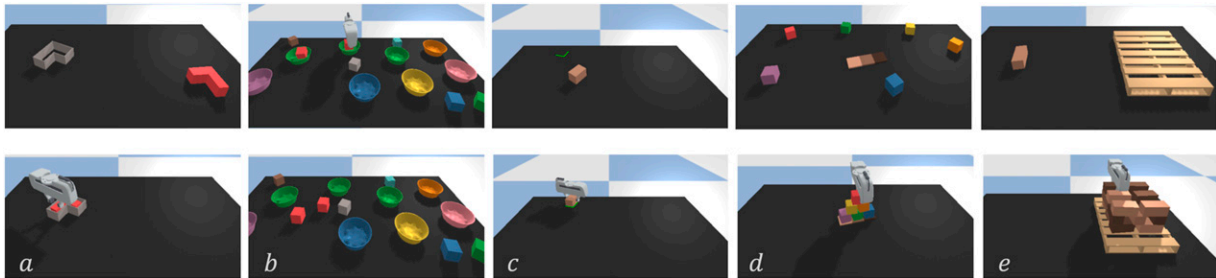
**6.3.3. Baselines.** We compare our method against Transporter Net ([Zeng et al., 2021](#)) as well as the following baselines previously used in the Transporter Net paper ([Zeng et al., 2021](#)). *Form2Fit* ([Zakka et al., 2020](#)) introduces a matching module with the measurement of  $L_2$  distance of high-dimension descriptors of picking and placing locations. *Conv-MLP* is a common end-to-end model ([Levine et al., 2016](#)) which outputs  $a_{\text{pick}}$  and  $a_{\text{place}}$  using convolution layers and MLPs (multi-layer perceptrons). *GT-State MLP* is a regression model composed of an MLP that accepts the ground-truth poses and 3D bounding boxes of objects in the environment. *GT-State MLP 2-step* outputs the actions sequentially with two MLP networks and feeds  $a_{\text{pick}}$  to the second step. All regression baselines learn mixture densities ([Bishop, 1994](#)) with log-likelihood loss.

For goal-conditioned tasks, we compare two baselines *Equivariant-Transporter-Goal-Stack* with *Transporter-Goal-Stack* and *Transporter-Goal-Split*.

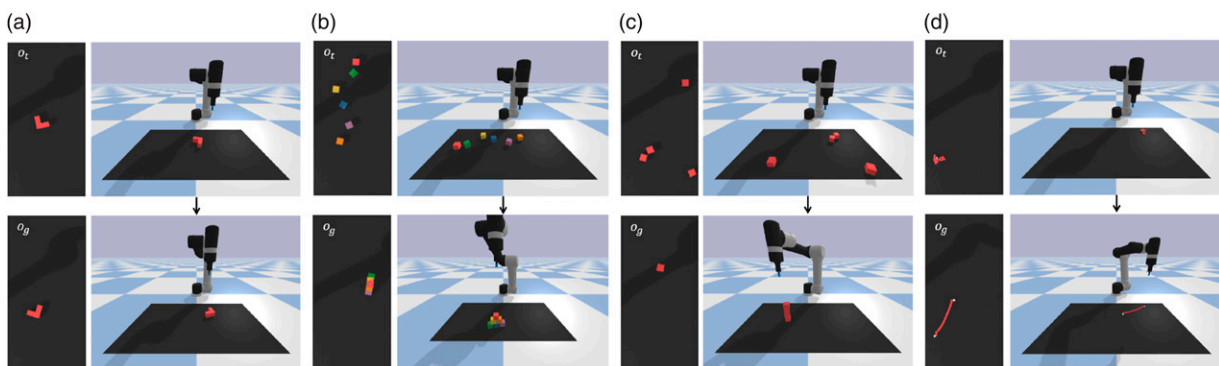
**6.3.4. Adaptation of transporter net for picking using a parallel-jaw gripper.** In order to compare our method against Transporter Net for the five parallel-jaw gripper tasks, we must modify Transporter to handle the gripper. We accomplish this by ([Zeng, Song, Welker, Lee, Rodriguez and Funkhouser, 2018a](#)) lifting the input scene image  $o_t$  over  $C_n$ , producing a stack of differently oriented input images which is provided as input to the pick network  $f_{\text{pick}}$ . The results are then counter-rotated at the output of  $f_{\text{pick}}$  and each channel corresponds to one pick orientation.

## 6.4. Results for the Ravens-10 benchmark tasks

**6.4.1. Task success rates.** [Table 1](#) shows the performance of our model on the Raven-10 tasks for different numbers of demonstrations used during training. All tests are evaluated



**Figure 7.** Simulated environment for parallel-jaw gripper tasks. From left to right: (a) inserting blocks into fixtures, (b) placing red boxes into green bowls, (c) align-box-corners to green lines, (d) stacking a pyramid of blocks, and (e) palletizing-boxes.



**Figure 8.** Simulated environment for goal-conditioned tasks. The first row shows the initial observations of four tasks where objects are generated with random poses on the workspace. The robot must use pick and place actions to manipulate the objects to the target pose specified in the goal images  $o_g$ , as shown in the second row. (a) Goal-conditioned block-insertion, (b) goal-conditioned block-pyramid, (c) goal-conditioned four blocks, and (d) goal-conditioned cable-align.

on unseen configurations, that is, random poses of objects, and three tasks (align-box-corner, assembling-kits, packing-box) use unseen objects. Our proposed Equivariant Transporter Net outperforms all the other baselines in nearly all cases. The amount by which our method outperforms others is largest when the number of demonstrations is smallest, that is, with only one or 10 demonstrations. With just 10 demonstrations per task, our method can achieve  $\geq 95\%$  success rate on 7/10 tasks. With either one or 10 demonstrations, the performance of our model is better than baselines trained with 1000 demonstrations on 5/10 tasks.

**6.4.2. Training efficiency.** Another interesting consequence of our more structured model is that training is much faster. Figure 9 shows task success rates as a function of the number of SGD steps for two tasks (Block-Insertion and Sweeping-Piles). Our equivariant model converges much faster in both cases. It indicates that the large symmetry group enables the model to learn on a low-dimension space and achieve better convergence speed.

### 6.5. Results for parallel-jaw gripper tasks

Table 2 compares the success rate of Equivariant Transporter with the baseline Transporter Net on parallel-jaw gripper tasks. Again, our method outperforms the

baseline in nearly all cases. One interesting observation that can be made by comparing Tables 1 and 2 is that both Equivariant Transporter and baseline Transporter do better on the gripper versions of the task compared to the original Ravens-10 versions. This is likely caused by the fact that the expert demonstrations we developed for the gripper version task have fewer stochastic gripper poses during pick than the case in the original Ravens-10 benchmark.

### 6.6. Results for goal-conditioned equivariant transporter net

Table 3 compares the performance of Equivariant-Transporter-Goal-Stack with the two baselines (Transporter-Goal-Stack, Transporter-Goal-Split) for goal-conditioned tasks. Our model gets better performance than the baselines on all the tasks. In most cases, the performance gap between our model and the baselines becomes larger as the number of demonstrations decreases. It shows the sample efficiency of our proposed method could be used to solve goal-conditioned tasks effectively.

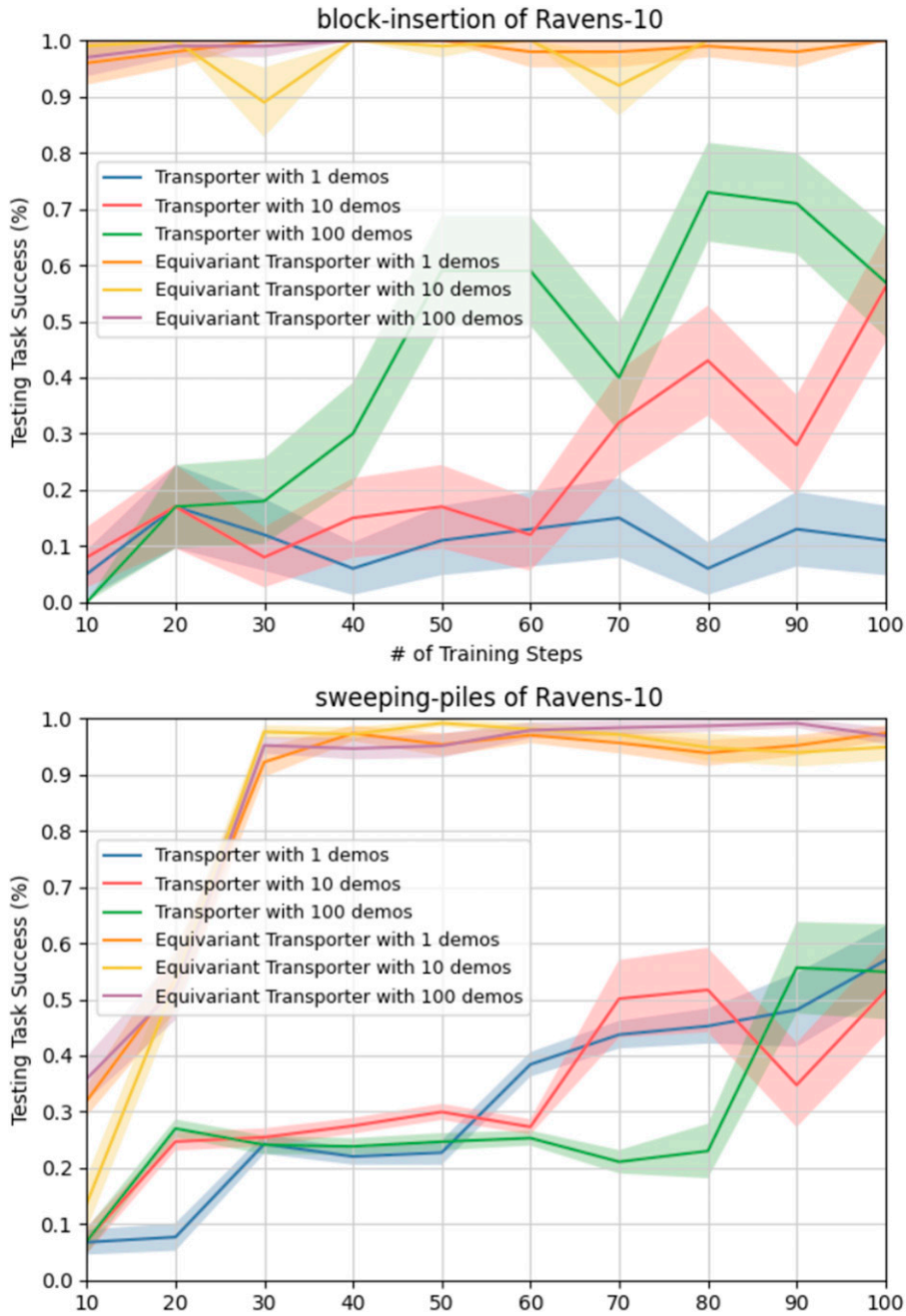
### 6.7. Ablation study

**6.7.1. Ablations.** We performed an ablation study to evaluate the relative importance of the equivariant models in

**Table 1.** Performance comparisons on ravens-10 benchmark (suction gripper). success rate (mean%) vs. the number of demonstration episodes (1, 10, 100, or 1000) used in training.

Method	Block-insertion				Place-red-in-green				Towers-of-Hanoi				Align-box-corner				Stack-block-pyramid			
	1	10	100	1000	1	10	100	1000	1	10	100	1000	1	10	100	1000	1	10	100	1000
Equivariant transporter	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>98.5</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>88.1</b>	<b>95.7</b>	<b>100</b>	<b>100</b>	41.0	<b>99.0</b>	<b>100</b>	<b>100</b>	<b>34.6</b>	<b>80.0</b>	<b>90.8</b>	<b>95.1</b>
Transporter network	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	84.5	<b>100</b>	<b>100</b>	<b>100</b>	73.1	83.9	97.3	98.1	35.0	85.0	97.0	98.0	13.3	42.6	56.2	78.2
Form2Fit	17.0	19.0	23.0	29.0	83.4	<b>100</b>	<b>100</b>	<b>100</b>	3.6	4.4	3.7	7.0	7.0	2.0	5.0	16.0	19.7	17.5	18.5	32.5
Conv. MLP	0.0	5.0	6.0	8.0	0.0	3.0	25.5	31.3	0.0	1.0	1.9	2.1	0.0	2.0	1.0	1.0	0.0	1.8	1.7	1.7
GT-state MLP	4.0	52.0	96.0	99.0	0.0	0.0	3.0	82.2	10.7	10.7	6.1	5.3	47.0	29.0	29.0	59.0	0.0	0.2	1.3	15.3
GT-state MLP 2-step	6.0	38.0	95.0	<b>100</b>	0.0	0.0	19.0	92.8	22.0	6.4	5.6	3.1	<b>49.0</b>	12.0	43.0	55.0	0.0	0.8	12.2	17.5
Palletizing-boxes																				
Assembling-kits					Packing-boxes					Manipulating-rope					Sweeping-piles					
1	10	100	1000	1	10	100	1000	1	10	100	1000	1	10	100	1000	1	10	100	1000	
Equivariant transporter	<b>75.3</b>	<b>98.9</b>	<b>99.6</b>	<b>99.6</b>	<b>63.8</b>	<b>90.6</b>	<b>98.6</b>	<b>100</b>	<b>98.3</b>	<b>99.4</b>	<b>99.6</b>	<b>100</b>	<b>31.0</b>	<b>85.0</b>	<b>92.3</b>	<b>98.4</b>	<b>97.9</b>	<b>99.5</b>	<b>100</b>	<b>100</b>
Transporter network	63.2	77.4	91.7	97.9	28.4	78.6	90.4	94.6	56.8	58.3	72.1	81.3	21.9	73.2	85.4	92.1	52.4	74.4	71.5	96.1
Form2Fit	21.6	42.0	52.1	65.3	3.4	7.6	24.2	37.6	29.9	52.5	62.3	66.8	11.9	38.8	36.7	47.7	13.2	15.6	26.7	38.4
Conv. MLP	31.4	37.4	34.6	32.0	0.0	0.2	0.2	0.0	0.3	9.5	12.6	16.1	3.7	6.6	3.8	10.8	28.2	48.4	44.9	45.1
GT-state MLP	0.6	6.4	30.2	30.1	0.0	0.0	1.2	11.8	7.1	1.4	33.6	56.0	5.5	11.5	43.6	47.4	7.2	20.6	63.2	74.4
GT-state MLP 2-step	0.6	9.6	32.8	37.5	0.0	0.0	1.6	4.4	4.0	3.5	43.4	57.1	6.0	8.2	41.5	58.7	9.7	21.4	66.2	73.9

Best performances are highlighted in bold.



**Figure 9.** Equivariant Transporter Network converges faster than Transporter Network. Top: Block-insertion task. Bottom: sweeping-piles task. On the block-insertion task, Equivariant Transporter can hit greater than 90% success rate after 10 training steps and achieve 100% success rate with less than 100 training steps.

**Table 2.** Performance comparisons on tasks with a parallel-jaw end effector. Success rate (mean%) vs. the number of demonstration episodes (1, 10, or 100) used in training..

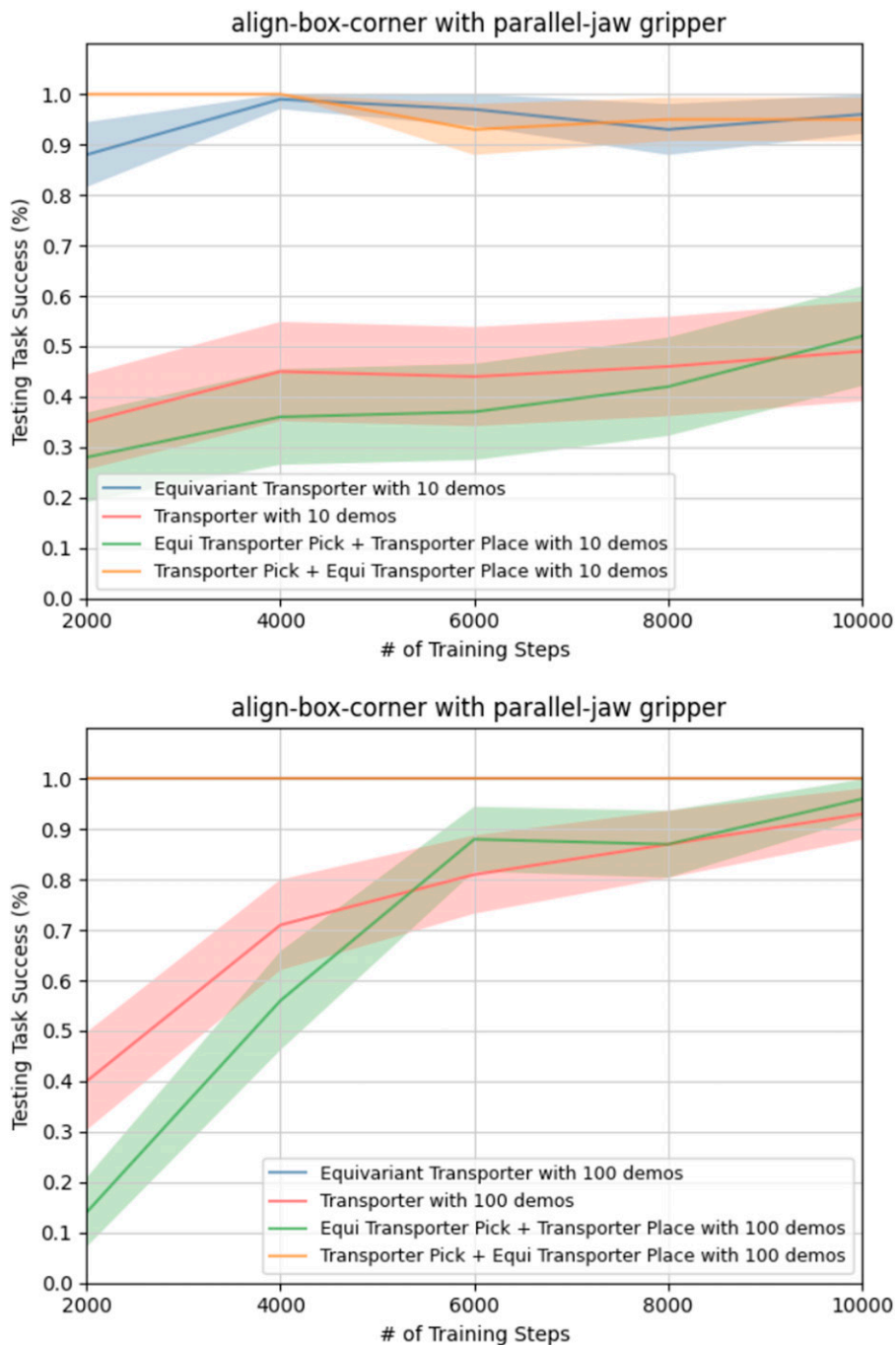
Method	Block-insertion			Place-red-in-green			Palletizing-boxes			Align-box-corner			Stack-block-pyramid		
	1	10	100	1	10	100	1	10	100	1	10	100	1	10	100
Equivariant transporter	<b>100</b>	<b>100</b>	<b>100</b>	<b>95.6</b>	<b>100</b>	<b>100</b>	<b>96.1</b>	<b>100</b>	<b>100</b>	<b>64.0</b>	<b>99.0</b>	<b>100</b>	<b>62.1</b>	<b>85.6</b>	<b>98.3</b>
Transporter network	98.0	<b>100</b>	<b>100</b>	82.3	94.8	<b>100</b>	84.2	99.6	<b>100</b>	45.0	85.0	99.0	16.6	63.3	75.0

Best performances are highlighted in bold.

**Table 3.** Performance comparisons on goal-conditioned tasks. Success rate (Mean%) vs. the number of demonstration episodes (1, 10, or 100) used in training.

Method	Goal-conditioned-block-insertion			Goal-conditioned-block-pyramid			Goal-conditioned-four-blocks			Goal-conditioned-cable-align		
	1	10	100	1	10	100	1	10	100	1	10	100
Equivariant-transporter-goal-stack	<b>100</b>	<b>100</b>	<b>100</b>	<b>51.3</b>	<b>84.7</b>	<b>86.5</b>	<b>63.5</b>	<b>87.3</b>	<b>93.6</b>	<b>74.9</b>	<b>89.4</b>	<b>92.6</b>
Transporter-goal-stack	<b>100</b>	<b>100</b>	<b>100</b>	9.8	64.7	72.5	11.2	24.1	25.3	49.5	78.7	88.7
Transporter-goal-split	99.0	<b>100</b>	<b>100</b>	3.3	58.8	67.7	17.8	27.3	27.9	52.5	84.4	92.1

Best performances are highlighted in bold.

**Figure 10.** Ablation study. Performance is evaluated on 100 unseen tests of each task.

pick ( $f_p$  and  $f_\theta$ ) and place ( $\psi$  and  $\phi$ ). We compare four versions of the model with various levels of equivariance: non-equivariant pick and non-equivariant place (baseline Transporter), equivariant pick and non-equivariant place, non-equivariant pick and equivariant place, and equivariant pick and equivariant place (Equivariant Transporter).

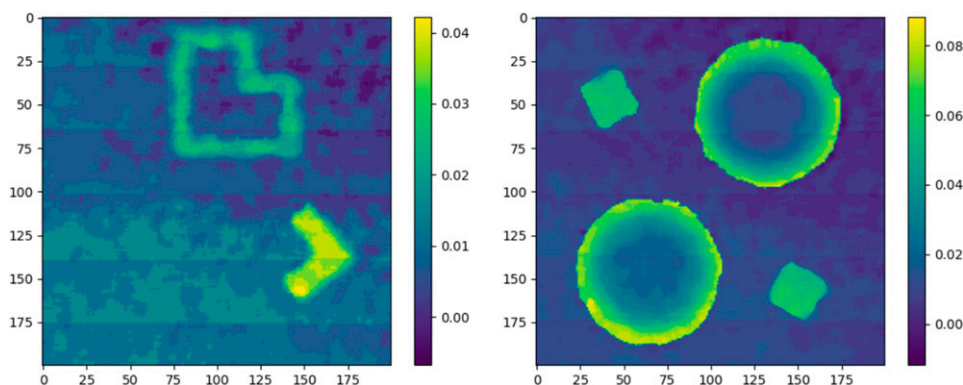
**6.7.2. Results.** Figure 10 shows the performance of all four ablations as a function of the number of SGD steps for the scenario where the agent is given 10 or 100 expert demonstrations. The results indicate that place equivariance (i.e., equivariance of  $\psi$  and  $\phi$ ) is, namely, responsible for the gains in performance of Equivariant Transporter versus baseline Transporter. This finding is consistent with the argument that the larger  $C_n \times C_n$  symmetry group (only present with the equivariant place) is responsible for our performance gains. Though the non-equivariant and equivariant pick networks result in comparable performance, the equivariant network is far more computationally efficient. The equivariant model takes a single image of the observation as input while the non-equivariant method (Zeng, Song, Welker, Lee, Rodriguez and Funkhouser, 2018a) needs a stack of  $n$ -different rotated input images in order to infer the pick orientation channel-wisely.

**6.7.3. Ablation study on group size.** We compare different group sizes to encode the latent features of our network. Note that the number of pick orientations and the number

of place orientations are task-relevant parameters and this ablation study is used to investigate the group size of the intermediate layers. We select  $(C_4, C_4, C_4)$ ,  $(C_6, C_6, C_6)$ ,  $(C_8, C_8, C_8)$  to construct three different settings of  $f_p$ ,  $\psi$  and  $\phi$ . We build a light version of our model with each setting and train it on block-insertion, packing-boxes, and manipulating-rope with 1 and 10 demos. Specifically, the  $C_4$  model has nine million trainable parameters, and the  $C_6$  model and the  $C_8$  model have 13.5 million and 18 million parameters, respectively. Note that a large group has a large fiber space dimension which results in more parameters but it also adds more constraints to the free parameters of the kernel. We train each model with data augmentation and evaluate the performance on 100 unseen test cases every 1k training steps. We report the highest success rate, its corresponding training steps, and the success rate with 1k training steps in Table 4. Several findings could be drawn from Table 4. First, the best mean success rates are consistent with different groups on most tasks. As the number of available demonstrations increases, the differences decrease. Second, large group size may boost the convergence speed when looking at the results of Block-Insertion-1. However, it could also result in overfitting when comparing the results of Manipulating-Rope-1 since the large-group model has more constraints and parameters. Finally, we think the group size of the intermediate layer might be regarded as a hyper-parameter.

**Table 4.** Ablation Study on Group Size of Intermediate Equivariant Layers. We report the success rate (% mean), and its corresponding training steps on three tasks from Ravens-10 Benchmark with 1 and 10 demos, respectively. We Test three different cyclic groups ( $C_4$ ,  $C_6$ ,  $C_8$ ) of the intermediate layers of the network  $f_p$ ,  $\psi$ , and  $\phi$ .

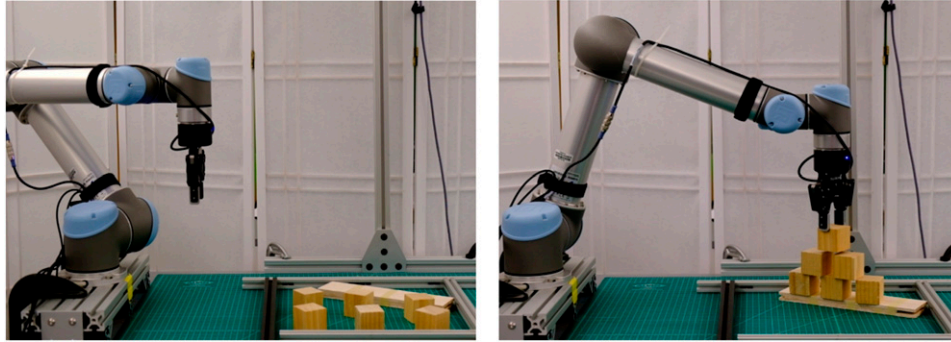
Task name-demo number group size	Block- insertion-1			Block- insertion-10			Packing- boxes-1			Packing- boxes-10			Manipulating- rope-1			Manipulating- Rope-10		
	$C_4$	$C_6$	$C_8$	$C_4$	$C_6$	$C_8$	$C_4$	$C_6$	$C_8$	$C_4$	$C_6$	$C_8$	$C_4$	$C_6$	$C_8$	$C_4$	$C_6$	$C_8$
Best mean success rate	100	100	100	100	100	100	98.3	99.6	98.2	99.4	99.8	99.4	47.0	50.7	40.2	81.0	83.9	79.9
SGD step of the best mean success rate	4k	1k	1K	1k	1k	1k	6k	5k	6k	2k	8k	6k	6k	7k	7k	1k	4k	5k
Mean success rate trained with 1k Step	92.0	100	100	100	100	100	87.7	94.1	97.2	99.2	97.2	98.8	28.0	39.0	12.3	81.0	70.9	43.2



**Figure 11.** Real robot experiment: initial observation  $o_t \in \mathbb{R}^{200 \times 200}$  from the depth sensor. The left figure shows the block-insertion task; the right figure shows the task of placing boxes in bowls. The depth value (meter) is illustrated in the color bar.

**Table 5.** Task success rates for physical robot evaluation tasks.

Task	# Demos	# Completions/# trials	Success rate(%)
Stack-block-pyramid	10	17/20	95.8
Place-box-in-bowl	10	20/20	100
Block-insertion	10	20/20	100

**Figure 12.** Stack-block-pyramid task on the real robot. The left figure shows the initial state; the right figure shows the completion state.

### 6.8. Experiments on a physical robot

We evaluated Equivariant Transporter on a physical robot in our lab. The simulator was not used in this experiment—all demonstrations were performed on the real robot.

**6.8.1. Setup.** We used a UR5 robot with a Robotiq-85 end effector. The workspace was a  $40\text{ cm} \times 40\text{ cm}$  region on a table beneath the robot (see Figure 12). The observations  $o$  were  $200 \times 200$  depth images obtained using an Occipital Structure Sensor that was mounted pointing directly down at the table (see Figure 11).

**6.8.2. Tasks.** We evaluated Equivariant Transporter on three of the Ravens-10 gripper-modified tasks: block-insertion, placing boxes in bowls, and stacking a pyramid of blocks. Since our sensor only measures depth (and not RGB), we modified the box-in-bowls task such that box color was irrelevant to success, that is, the task is simply to put any box into a bowl.

**6.8.3. Demonstrations.** We obtained 10 human demonstrations of each task. These demonstrations were obtained by releasing the UR5 brakes and pushing the arm physically so that the harmonic actuators were back-driven.

**6.8.4. Training and evaluation.** For each task, a single-policy agent was trained for 10k SGD steps. During testing, objects were randomly placed on the table. A task was considered to have failed when a single incorrect pick or place occurred. We evaluated the model on 20 unseen configurations of each task.

**6.8.5. Results.** Table 5 shows results from 20 runs of each of the three tasks. Notice that the success rates here are

higher than they were for the corresponding tasks performed in simulation (Table 2). This is likely caused by the fact that the criteria for task success in simulation (less than 1 cm translational error and less than  $\pi/12$  rotation error) were more conservative than is actually the case in the real world.

### 6.9. Discussion

Equivariant networks are built on top of conventional convolution kernels with the steerability constraint. It does not break the mechanism of weight sharing and updating and thus keeps the robustness of learning and reasoning of CNN. As shown in Figure 11, Equivariant Transporter Net can handle real-sensor noise. Frequently, the crop  $c$  contains multiple objects. For instance, on the stack-block-pyramid task as shown in Figure 12, the crop not only includes the block to be picked but also neighboring blocks or some parts of them. During training, data augmentation also generates images with partially observed objects. For example, on the block-insertion task, it shifts some part of the L-shape block or the slot out of the scene. Some special shapes like elongated objects might be difficult to represent with an image crop and may be suitable for the goal-conditioned version of our model. Some high-precision tasks like gear assembly are more sensitive to discretization and it may be tackled more easily with the  $SO(2)$  version of our model.

Compared to learning pick and place skills efficiently, the one-shot generalization and sequential decision-making ability of both Transporter Net and Equivariant Transporter Net seem less compelling. As shown in Table 1, they achieved less than 50% success rate when trained with one demo on the align-box-corner task that requires the agent to generalize the skill to boxes with random sizes and colors during the test. The performances on the stack-block-pyramid task trained with one demo are below 40%

since if there was a collapse, some blocks might be tilted and it yields out-of-distribution data.

## 7. Conclusion and limitations

This paper explores the symmetries present in the pick and place problem and finds that they can be described by pick symmetry and place symmetry. This corresponds to the group of different pick and place poses. We evaluate the Transporter Network model proposed in Zeng et al. (2021) and find that it encodes half of the place symmetry (the  $C_n$ -place symmetry). We propose a novel version of Transporter Net, Equivariant Transporter Net, which we show encodes both types of symmetries. The large symmetry group could also be extended to solve goal-conditioned tasks. We evaluate our model on the Ravens-10 Benchmark and its variations and evaluate against multiple strong baselines. Finally, we demonstrate that the method can effectively be used to learn manipulation policies on a physical robot.

One limitation of our framework as it is presented in this paper is that it relies entirely on behavior cloning. A clear direction is to integrate more on-policy learning which we believe would enable us to handle more complex tasks. Other directions of the multi-task language-conditioned equivariant agent, a closed-loop policy, or 3D Equivariant Transporter Net are also interesting.

### Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

### Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This work was supported in part by NSF 1724257, NSF 1724191, NSF1763878, NSF 1750649, NSF 2107256, NSF 2134178, NASA 80NSSC19K1474, the Harold Alfond Foundation, and the Roux Institute.

### ORCID iD

Haojie Huang  <https://orcid.org/0000-0001-8737-7959>

### Notes

1. Our implementation uses the discrete  $C_n$  group instead of the continuous  $SO(2)$  group in order to compare with the baseline *Transporter Net* (Zeng et al., 2021). The  $SO(2)$  version of our model could be easily achieved with the irreducible representations based on our implementation.
2. The  $SO(2)$  equivariance can be approximately achieved with irreducible representations as the output type for  $f_{\theta}$ .
3. The  $SO(2)/C_2$  quotient group can be realized by using the basis function with a period of  $\pi$ .
4. The resolution of the input image is  $320 \times 160$  for this training time.

## References

- Berscheid L, Meisner P and Kröger T (2020) Self-supervised learning for precise pick-and-place without object model. *IEEE Robotics and Automation Letters* 5(3): 4828–4835.
- Besl PJ and McKay ND (1992) Method for registration of 3-d shapes. In: *Sensor fusion IV: Control Paradigms and Data Structures*. Bellingham, Washington: International Society for Optics and Photonics, Vol. volume 1611, 586–606.
- Bishop CM (1994) Mixture density networks.
- Cesa G, Lang L and Weiler M (2021) A program to build e (n)-equivariant steerable cnns. In: *International Conference on Learning Representations*. Vienna, Austria: ICLR
- Chen X, Chen R, Sui Z, et al. (2019) Grip: generative robust inference and perception for semantic robot manipulation in adversarial environments. In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Macau, China: IEEE, 3988–3995.
- Chen H, Liu S, Chen W, et al. (2021) Equivariant point network for 3d point cloud analysis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. New York, NY, USA: IEEE, 14514–14523.
- Cohen T and Welling M (2016) Group equivariant convolutional networks. In: *International Conference on Machine Learning*. Vienna, Austria: PMLR, 2990–2999.
- Cohen TS and Welling M (2017) Steerable cnns. In: *International Conference on Learning Representations*. Vienna, Austria: ICLR.
- Curtis A, Fang X, Kaelbling LP, et al. (2022) Long-horizon manipulation of unknown objects via task and motion planning with estimated affordances. In: *2022 International Conference on Robotics and Automation (ICRA)*. Philadelphia, PA, USA: IEEE, 1940–1946.
- Deng X, Xiang Y, Mousavian A, et al. (2020) Self-supervised 6d object pose estimation for robot manipulation. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. Paris, France: IEEE, 3665–3671.
- Deng C, Litany O, Duan Y, et al. (2021) Vector neurons: a general framework for so (3)-equivariant networks. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. Montreal, BC, Canada: IEEE, 12200–12209.
- Devin C, Rowghanian P, Vigorito C, et al. (2020) Self-supervised goal-conditioned pick and place. arXiv preprint arXiv: 2008.11466.
- Fuchs F, Worrall D, Fischer V, et al. (2020) Se (3)-transformers: 3d roto-translation equivariant attention networks. *Advances in Neural Information Processing Systems* 33: 1970–1981.
- Gualtieri M and Platt R (2021) Robotic pick-and-place with uncertain object instance segmentation and shape completion. *IEEE Robotics and Automation Letters* 6(2): 1753–1760.
- He K, Zhang X, Ren S, et al. (2016) Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. New York, NY, USA: IEEE, 770–778.
- Hester T, Vecerik M, Pietquin O, et al. (2018) Deep q-learning from demonstrations. In: *Thirty-second AAAI Conference on Artificial Intelligence*. Washington, DC, US: AAAI.

- Huang H, Yang Z and Platt R (2021) Gascn: graph attention shape completion network. In: *2021 International Conference on 3D Vision (3DV)*. London, UK: IEEE, 1269–1278.
- Huang H, Wang D, Walters R, et al. (2022a) Equivariant transporter network. In: *Proceedings of Robotics: Science and Systems*. New York, NY, USA: IEEE. DOI: [10.15607/RSS.2022.XVIII.007](https://doi.org/10.15607/RSS.2022.XVIII.007).
- Huang H, Wang D, Zhu X, et al. (2022b) Edge grasp network: a graph-based se (3)-invariant approach to grasp detection. arXiv preprint arXiv:2211.00191.
- Hussein A, Gaber MM, Elyan E, et al. (2017) Imitation learning: a survey of learning methods. *ACM Computing Surveys* 50(2): 1–35.
- Jaegle A, Borgeaud S, Alayrac JB, et al. (2021) Perceiver Io: a general architecture for structured inputs and outputs. arXiv preprint arXiv:2107.14795.
- Jenner E and Weiler M (2021) Steerable partial differential operators for equivariant neural networks. arXiv preprint arXiv:2106.10163.
- Jia M, Wang D, Su G, et al. (2023) Seil: simulation-augmented equivariant imitation learning. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. ExCeL London: IEEE, 1845–1851.
- Khansari M, Kappler D, Luo J, et al. (2020) Action image representation: learning scalable deep grasping policies with zero real world data. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. Paris, France: IEEE, 3597–3603.
- Kingma DP and Ba J (2014) Adam: a method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Lang L and Weiler M (2020) A wigner-eckart theorem for group equivariant convolution kernels. arXiv preprint arXiv:2010.10952.
- Levine S, Finn C, Darrell T, et al. (2016) End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research* 17(1): 1334–1373.
- Liu X, Jonschkowski R, Angelova A, et al. (2020) Keypose: multi-view 3d labeling and keypoint estimation for transparent objects. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. New York, NY, USA: IEEE, 11602–11610.
- Manuelli L, Gao W, Florence P, et al. (2019) kpm: keypoint affordances for category-level robotic manipulation. In: *The International Symposium of Robotics Research*. Berlin, Germany: Springer, 132–157.
- Morrison D, Leitner J and Corke P (2018) Closing the loop for robotic grasping: a real-time, generative grasp synthesis approach. In: *Proceedings of Robotics: Science and Systems*. Pittsburgh, PA, USA: Carnegie Mellon University. DOI: [10.15607/RSS.2018.XIV.021](https://doi.org/10.15607/RSS.2018.XIV.021).
- Nagabandi A, Konolige K, Levine S, et al. (2020) Deep dynamics models for learning dexterous manipulation. In: *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, 1101–1112.
- Nair V and Hinton GE (2010) Rectified linear units improve restricted Boltzmann machines In: *Icml*. Vienna, Austria: International Conference on Machine Learning
- Narayanan V and Likhachev M (2016) Discriminatively-guided deliberative perception for pose estimation of multiple 3d object instances. In: *Robotics: Science and Systems*. Delft, Netherlands: Delft University of Technology.
- Qureshi A, Mousavian A, Paxton C, et al. (2021) Nerp: neural rearrangement planning for unknown objects. In: *Proceedings of Robotics: Science and Systems (RSS)*. Pittsburgh, PA, USA: Carnegie Mellon University.
- Ronneberger O, Fischer P and Brox T (2015) U-net: convolutional networks for biomedical image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Berlin, Germany: Springer, 234–241.
- Seita D, Florence P, Tompson J, et al. (2021) Learning to rearrange deformable cables, fabrics, and bags with goal-conditioned transporter networks. In: *IEEE International Conference on Robotics and Automation (ICRA)*. Xi'an, China: IEEE.
- Serre JP (1977) *Linear representations of finite groups*. Berlin, Germany: Springer, vol. 42.
- Shridhar M, Manuelli L and Fox D (2022a) Cliport: what and where pathways for robotic manipulation. In: *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, 894–906.
- Shridhar M, Manuelli L and Fox D (2022b) Perceiver-actor: a multi-task transformer for robotic manipulation. In: *Proceedings of the 6th Conference on Robot Learning (CoRL)*. Auckland, New Zealand: PMLR.
- Simeonov A, Du Y, Tagliasacchi A, et al. (2022) Neural descriptor fields: Se (3)-equivariant object representations for manipulation. In: *2022 International Conference on Robotics and Automation (ICRA)*. Philadelphia, PA, USA: IEEE, 6394–6400.
- Thomas N, Smidt T, Kearnes S, et al. (2018) Tensor field networks: rotation-and translation-equivariant neural networks for 3d point clouds. arXiv preprint arXiv:1802.08219.
- Vecerik M, Hester T, Scholz J, et al. (2017) Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. arXiv preprint arXiv:1707.08817.
- Wang D, Kohler C and Platt R (2021) Policy learning in se (3) action spaces. In: *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, 1481–1497.
- Wang D, Walters R, Zhu X, et al. (2022) Equivariant  $q$  learning in spatial action spaces. In: *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, 1713–1723.
- Weiler M and Cesa G (2019) General e (2)-equivariant steerable cnns. *Advances in Neural Information Processing Systems* 32: 14357–14368.
- Weiler M, Geiger M, Welling M, et al. (2018) 3d steerable cnns: learning rotationally equivariant features in volumetric data. *Advances in Neural Information Processing Systems* 31: 10402–10413.
- Wen B, Lian W, Bekris K, et al. (2022) You only demonstrate once: category-level manipulation from single visual demonstration. In: *Proceedings of Robotics: Science and Systems*. New York City, NY, USA: Columbia University. DOI: [10.15607/RSS.2022.XVIII.044](https://doi.org/10.15607/RSS.2022.XVIII.044).
- Wu Y, Yan W, Kurutach T, et al. (2020) Learning to manipulate deformable objects without demonstrations. In: *16th*

- Robotics: Science and Systems, RSS 2020*. Cambridge, MA, USA; MIT Press Journals.
- Xue Z, Yuan Z, Wang J, et al. (2022) Useek: unsupervised Se (3)-equivariant 3d keypoints for generalizable manipulation. arXiv preprint arXiv:2209.13864.
- Yoon Y, DeSouza GN and Kak AC (2003) Real-time tracking and pose estimation for industrial objects using geometric features. In: *2003 IEEE International Conference on Robotics and Automation (Cat. No. 03CH37422)*. Taipei, Taiwan: IEEE, volume 3, pp. 3473–3478.
- Yuan W, Khot T, Held D, et al. (2018) Pcn: point completion network. In: *2018 International Conference on 3D Vision (3DV)*. Verona, Italy: IEEE, pp. 728–737.
- Zakka K, Zeng A, Lee J, et al. (2020) Form2fit: learning shape priors for generalizable assembly from disassembly. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. Paris, France: IEEE, 9404–9410.
- Zeng A, Song S, Welker S, et al. (2018a) Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Madrid, Spain: IEEE, pp. 4238–4245.
- Zeng A, Song S, Yu KT, et al. (2018b) Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching. In: *2018 IEEE International Conference on Robotics and Automation (ICRA)*. Brisbane, QLD, Australia: IEEE, 3750–3757.
- Zeng A, Florence P, Tompson J, et al. (2021) Transporter networks: rearranging the visual world for robotic manipulation. In: *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, 726–747.
- Zhu X, Wang D, Biza O, et al. (2022) Sample efficient grasp learning using equivariant models. In: *Proceedings of Robotics: Science and Systems (RSS)*. New York, NY, USA: Columbia University.