

# One-Shot Imitation Learning with Invariance Matching for Robotic Manipulation

Xinyu Zhang and Abdeslam Boularias  
Rutgers University  
Email: {xz653, ab1544}@rutgers.edu

**Abstract**—Learning a single universal policy that can perform a diverse set of manipulation tasks is a promising new direction in robotics. However, existing techniques are limited to learning policies that can only perform tasks that are encountered during training, and require a large number of demonstrations to learn new tasks. Humans, on the other hand, often can learn a new task from a single unannotated demonstration. In this work, we propose the Invariance-Matching One-shot Policy Learning (IMOP) algorithm. In contrast to the standard practice of learning the end-effector’s pose directly, IMOP first learns invariant regions of the state space for a given task, and then computes the end-effector’s pose through matching the invariant regions between demonstrations and test scenes. Trained on the 18 *RLBench* tasks, IMOP achieves a success rate that outperforms the state-of-the-art consistently, by 4.5% on average over the 18 tasks. More importantly, IMOP can learn a novel task from a single unannotated demonstration, and without any fine-tuning, and achieves an average success rate improvement of 11.5% over the state-of-the-art on 22 novel tasks selected across nine categories. IMOP can generalize to new shapes and objects that are different from those in the demonstration. Further, IMOP can perform one-shot sim-to-real transfer using a single real-robot demonstration.

## I. INTRODUCTION

Multi-Task Learning (MTL) is a recent type of learning technique where a single policy is trained to perform various tasks. With the recent advances in computer vision architectures [20], a popular MTL strategy in robotics is to acquire a multi-task control policy through imitation learning from visual demonstrations [16, 17, 38, 41]. These methods achieved impressive results in performing challenging manipulation tasks in unstructured 3D environments, such as *RLBench* [24]. However, these methods rely on the assumption that training and testing share the same set of tasks. Moreover, training such policies on new tasks requires hundreds of demonstrations [17], which results in the catastrophic forgetting of old tasks [26].

To address these issues, one-shot imitation learning aims to learn on a set of *base tasks* and generalize to *novel tasks* given only a single demonstration for each novel task and without re-training. However, existing methods rely on strong assumptions, such as requiring the novel tasks to be limited variations of the same base tasks [42], requiring the base and novel tasks to have the same object setup [11], only generalizing specific actions on certain categories of objects with known 3D models [2], or operating in simple 2D planar environments [13]. Moreover, most recent works focus on applying general popular

techniques, e.g., transformers and diffusion models, to one-shot imitation settings, without taking advantage of the particular structure of robotic manipulation tasks [11, 32, 42, 43].

Therefore, a key question here is: can robots learn a manipulation policy that not only performs well on base tasks but also generalizes to novel unseen tasks using a single demonstration and without any fine-tuning? To this end, we propose IMOP (Invariance Matching One-shot Policy Learning), a new algorithm that not only outperforms the state-of-the-art on the standard 18 tasks of *RLBench* (69.6% average overall success rate compared to 65.1% for the state-of-the-art [16]), but also generalizes to 22 novel tasks with a single demonstration and without any fine-tuning (41.3% average success rate compared to 29.8% for the state-of-the-art). The 22 novel tasks are selected across nine categories and are substantially different from the base tasks in terms of objects and task goals [19, 24]. Moreover, we find that IMOP also generalizes to substantial shape variations, and can manipulate new objects that are different from those in the demonstration.

Instead of learning the desired end-effector’s pose directly, IMOP learns key *invariant regions* of each task, and finds pairwise correspondences between the invariant regions in a one-shot demonstration and in a given test scene. The pairwise correspondences are used to analytically compute the desired end-effector’s pose in the test scene from the least-squares solution of a point-set registration problem. The invariant region is defined as a set of 3D points whose coordinates remain invariant when viewed in the end-effector’s frame, across states that share the same semantic action. We devise a graph-based invariant region matching network. The invariant regions are located through neighbor attention [48] from the KNN graphs that connect point clouds of demonstration and test scenes. The ground-truth invariant regions are discovered in offline unannotated demonstrations of base tasks.

To summarize, our contributions are threefold. (1) We propose IMOP, a one-shot imitation learning algorithm for robotic manipulation that learns a universal policy that is not only successful on base tasks, but also generalizes to novel tasks using a single unannotated demonstration. (2) We propose a correspondence-based pose regression method for manipulation tasks, which predicts robot actions by matching key visual elements, and a graph-based invariant region matching network on KNN graphs that connect demonstrations and test scenes. (3) We present a thorough empirical study of the performance and generalization ability of IMOP on a diverse set of tasks.

## II. RELATED WORK

### Learning robotic manipulation from demonstrations.

Learning manipulation policies from offline visual demonstrations has gained increasing attention following the rapid improvement of vision models [12, 28, 29, 3, 4, 36]. Prior works such as Transporter networks [45] and CLIPort [37] learned simple pick-and-place tasks in a 2D planar setting. C2F-ARM [25] and PerAct [38] extend the control capacity to the 3D space with 6-DoF but learn a separate policy for each single task. More recently, RVT [17], Hiveformer [19], Act3D [16], and ChainedDiffuser [41] learned a multi-task policy from demonstrations of multiple tasks. Yet, these prior approaches can only perform tasks that were seen during training and cannot generalize to novel tasks upon inference. Training these models also requires hundreds of demonstrations per task. Instead, this work aims to develop a multi-task policy that not only performs well on seen tasks but also generalizes to novel tasks given a single demonstration.

**One-shot imitation learning.** Traditional imitation learning considers learning a policy for a single task with many expert trajectories provided for this particular task. One-shot imitation learning aims to learn from demonstrations of a set of base tasks and generalize to novel tasks using a single trajectory per new task without further training. Duan et al. [13] trained a one-shot imitator that generalizes to simple task variations in a 2D stacking environment. Finn et al. [14] used the same setting in the context of meta-learning. Dasari and Gupta [11], Xu et al. [42] and Mandi et al. [32] trained transformer-based one-shot imitators, such as encoding expert trajectories as prompts. Biza et al. [2] achieved one-shot generalization for grasping tasks where 3D models of the objects are available. Xu et al. [43] trained adaptation layers through hyper-networks. Yet, none of these prior works considered 6D manipulation tasks in the context of one-shot learning. Moreover, recent works focus on applying existing techniques, such as transformers and adaptation layers, to one-shot imitation settings, but did not investigate the inherent task structure of robotic manipulation. In contrast, this work proposes a one-shot imitation learner for 6D manipulation by discovering invariant regions of each task.

**Invariance and affordance in manipulation.** Incorporating invariance into deep learning models has been shown to drastically increase data efficiency and generalization [5, 7, 8, 31, 39]. Graf et al. [18] design image augmentations to learn viewpoint-invariant features for manipulation. Visual affordance is a region representation of action possibility [1, 34]. For example, Where2Act learns a probability map for pushing and pulling at each pixel [35]. In this work, we propose the concept of *invariant region*. Instead of camera viewpoint invariance, we train neural networks to predict regions whose positions remain invariant to the robot end-effector for a given task. Unlike affordance, the proposed invariant region is not used to represent action probability but to transfer actions from demonstrations to test scenes.

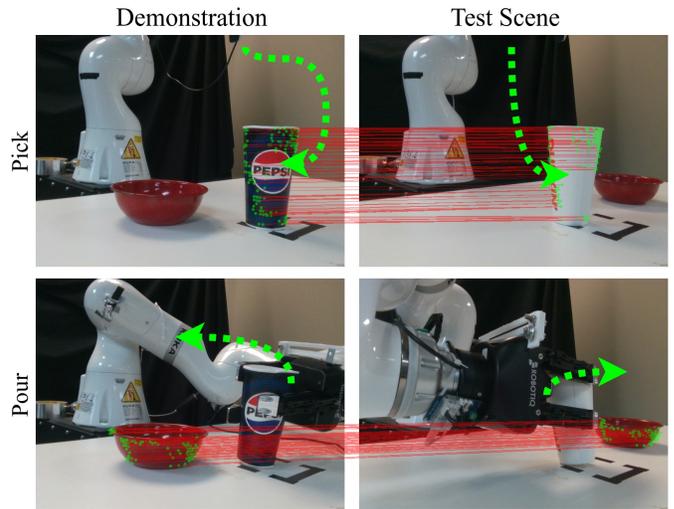


Fig. 1: Example of a pick-and-pour task executed by a real Kuka robot after observing a single demonstration and without any re-training. The correspondences between estimated invariant regions of the 3D point cloud in the demonstration and in the test scene are visualized in red lines. The invariant regions are predicted and matched by a neural network, trained offline.

## III. ONE-SHOT IMITATION WITH INVARIANCE MATCHING

### A. A Motivating Example

Figure 1 illustrates an example from our experiments where a Kuka robot is tasked with picking up a cup and pouring it into a bowl, using a single demonstration of picking up and pouring a different cup in a different position. At the core of IMOP lies the capacity to estimate and match invariant regions of the given task. By finding correspondences of invariant regions in the demonstration and in the test scene, the demonstrated actions can be transferred to the test scene. In this example, the invariant regions are: (1) a set of 3D points on the surfaces of the cups where contact with the fingertips occurs, and (2), a set of 3D points from the bowl’s point cloud that capture its spherical concave shape. The invariant regions of the 3D point cloud are learned offline from unannotated demonstrations, on various objects and tasks.

### B. Formulation of One-shot Manipulation Learning Problem

We use  $\mathbf{T}$  to denote a set of manipulation tasks,  $\mathbf{T} = \{\text{Task}_1, \text{Task}_2, \dots, \text{Task}_{N_{\text{tasks}}}\}$ . Each manipulation task  $\text{Task}_i \in \mathbf{T}$  is performed in a Markov Decision Process (MDP) represented by 5-tuple  $\mathcal{M}_i = (\mathcal{S}, \mathcal{A}, \mathbf{P}, R_i, \mu_i)$ , wherein  $\mathcal{S}$  and  $\mathcal{A}$  are state and action spaces,  $\mathbf{P} : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  is a state transition probability,  $R_i : \mathcal{S} \mapsto \mathbb{R}$  denotes a reward function associated with task  $\text{Task}_i$ , and  $\mu_i$  is the initial state distribution of  $\text{Task}_i$ .  $\mathbf{T}$  is composed of a set of base tasks, denoted by  $\mathbf{T}^{\text{base}}$ , and a set of novel tasks, denoted by  $\mathbf{T}^{\text{novel}}$ . Thus,  $\mathbf{T} = \mathbf{T}^{\text{base}} \cup \mathbf{T}^{\text{novel}}$ . During training, a large collection of offline demonstrations is provided for the base tasks. The underlying reward function of each task is unknown and is not used in this work. During testing, only one successful

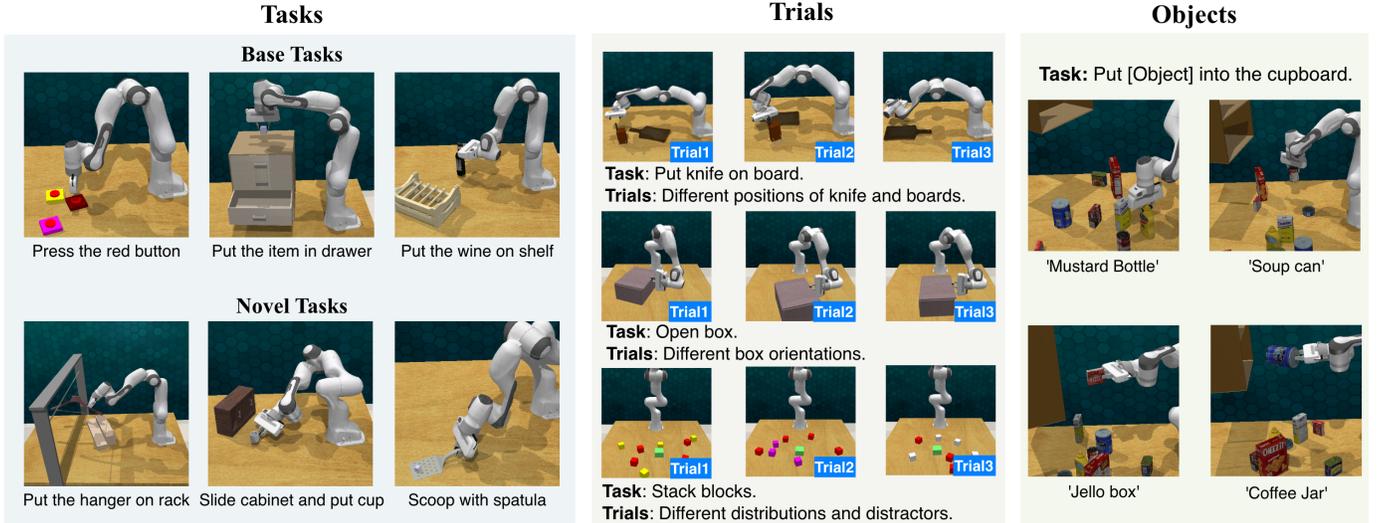


Fig. 2: Examples of tasks, and object-level generalization in various trials. After training on base tasks, IMOP is evaluated on novel tasks that are substantially different from the base tasks. Every learned task is evaluated in multiple trials, each with different object layouts and orientations. For each novel task, only one recorded trajectory is given as a demonstration. We summarize the performance of IMOP on both base and novel tasks in Section IV-A and IV-B. Furthermore, we study the ability of IMOP to generalize to new objects beyond those in the demonstrations in Section IV-C.

trajectory  $\tau_i = \{(s_i, a_i, s'_i)\}^{|\tau_i|}$  is provided to the robot for a novel Task $_i \in \mathbf{T}^{novel}$ , where  $|\tau_i|$  denotes the number of transitions,  $a_i$  is the demonstrated action in state  $s_i$ , and  $s'_i$  is the resulting next state. The goal is to learn a one-shot universal policy  $\pi(s, \tau_i)$ , which is a multi-task policy that returns an action  $a$  in new state  $s$  given a single demonstrated trajectory  $\tau_i$  of the desired novel skill.  $\tau_i$  is provided in an environment that is different from the one in which  $\pi$  is tested.

A state  $s \in \mathcal{S}$  is a raw 3D point cloud of the entire scene, including the robot. Each point in state  $s$  is represented in homogeneous coordinates. An action  $a \in \mathcal{A}$  is the desired (i.e., target) 6D pose of the robot’s end-effector, along with the desired state of the gripper. To execute an action, low-level robot movements are generated using off-the-shelf motion planners. More precisely, we define actions as 18-dimensional vectors  $a = (T, \lambda)$ , wherein  $T \in \mathbb{R}^{4 \times 4}$  denotes the target end-effector pose as a homogeneous matrix and is referred to as action pose, and  $\lambda$  denotes two binary values that indicate whether gripper is open and collisions are allowed or not.

### C. Proposed Method

1) *Invariant Region Matching Network*: We first consider the following problem. Given a task Task $_k \in \mathbf{T}$ , an observed transition  $(s_i, a_i, s'_i)$  in a demonstrated trajectory, and a new state  $s_j$  such that  $s_i \equiv s_j$ , how can action pose  $T_i$  in state  $s_i$  be translated to a new action pose  $T_j$  to execute in the new state  $s_j$ ? We use  $s_i \equiv s_j$  to indicate that  $s_i$  and  $s_j$  share the same optimal manipulation action. For example, the first manipulation action of the task ‘open door’ is ‘reach the door handle’. Therefore, states that describe different scenes with closed doors all share the optimal action ‘reach the door handle’. We refer to  $s_i$  as the support state and to the new state  $s_j$  as

the query state.

To answer this question, we define the notion of **invariant regions** in the following.

**Definition III.1.** The invariant region of state  $s_i$  in Task $_k$  is defined as  $\mathcal{I}(s_i | \text{Task}_k) = \{p \in s_i | \forall s_j \in \mathcal{S} \text{ s.t. } s_i \equiv s_j, \exists q \in s_j : \|T_i^{-1}p - T_j^{-1}q\| < \epsilon\}$ , where  $(T_i, \lambda_i) = \pi^*(s_i | \text{Task}_k)$  and  $(T_j, \lambda_j) = \pi^*(s_j | \text{Task}_k)$ .

In this definition,  $s_i$  is a scene point-cloud, and  $\mathcal{I}(s_i | \text{Task}_k)$  is the set of all the 3D points of  $s_i$  whose coordinates in frame  $T_i$  remain invariant across all the scenes  $s_j$  that share the same optimal action as the optimal action of  $s_i$ . Frame  $T_i$  is the frame of the end-effector’s pre-contact pose when performing the optimal manipulation action in scene  $s_i$ . For example, for the task of picking up a cup, the optimal action is to grasp the cup’s handle, regardless of the location, orientation, or style of the cup. The learned invariant region will likely be the cup’s handle. To simplify the notation, we drop Task $_k$  and simply refer to the invariant region as  $\mathcal{I}(s_i)$ , while it is implied that  $\mathcal{I}(s_i)$  is the invariant region of  $s_i$  for a given Task $_k$ .

The optimal action pose  $T_j$  of the query scene  $s_j$  can be computed by first predicting the invariant regions  $\mathcal{I}(s_i)$  and  $\mathcal{I}(s_j)$  of  $s_i$  and of support scene  $s_j$ , and then transforming  $T_i$  according to a matching between  $\mathcal{I}(s_i)$  and  $\mathcal{I}(s_j)$ . The overall framework of our proposed invariant region prediction and matching network is shown in Figure 3. We first build a KNN graph for each scene point-cloud by connecting each point to its  $k$  nearest points within the same scene. Next, we apply graph self-attention within each support scene  $s_i$ , and cross-attention between the KNN graphs of pairs of consecutive frames  $s_i$  and  $s'_i$  within the same support demonstration. In contrast to traditional attention that is computed globally over all points,

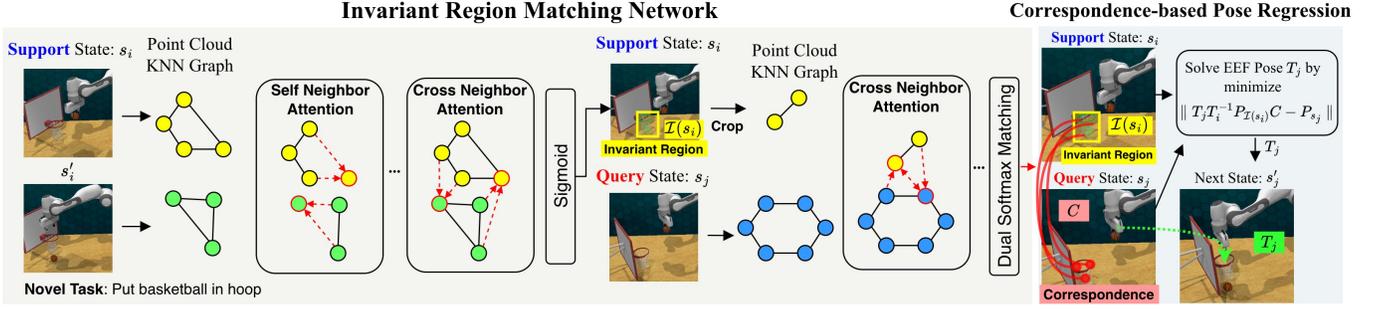


Fig. 3: Overview of the proposed invariant region matching network and correspondence-based pose regression. Given the observed transition of the support state  $s_i$ , we build KNN graphs from the scene point cloud and apply graph self-attention (within  $s_i$  and  $s'_i$ ) and cross-attention (between  $s_i$  and  $s'_i$ ) layers. The invariant region  $\mathcal{I}(s_i)$  is predicted as the set of activated points through a point-wise sigmoid over the support state  $s_i$ . Next, we apply graph cross-attention layers between the KNN graphs of  $\mathcal{I}(s_i)$  and query state  $s_j$ . The correspondence matrix  $C$  is predicted by matching the point-wise features  $h_{\mathcal{I}(s_i)}$  and  $h_{s_j}$ . The action pose  $T_j$  of query state  $s_j$  is analytically solved from the action pose  $T_i$  of support state  $s_i$ , correspondence matrix  $C$ ,  $P_{\mathcal{I}(s_i)}$  and  $P_{s_j}$  (points in  $\mathcal{I}(s_i)$  and  $s_j$ ), as detailed in Section III-C2.

graph attention operates within the local neighborhood of each given point. We use the point transformer layer [40] as the graph attention operator. The invariant region  $\mathcal{I}(s_i)$  is predicted as the set of activated points through a point-wise sigmoid over  $s_i$ . Then, we apply graph cross-attention layers between the KNN graph of  $\mathcal{I}(s_i)$  and the query state  $s_j$  to extract the point-wise features  $h_{\mathcal{I}(s_i)} \in \mathbb{R}^{|\mathcal{I}(s_i)| \times D}$  and  $h_{s_j} \in \mathbb{R}^{|s_j| \times D}$ , where  $D$  denotes the size of feature dimension. Finally, we perform a dual softmax matching [27] between  $h_{\mathcal{I}(s_i)}$  and  $h_{s_j}$  to obtain the correspondence matrix  $C \in [0, 1]^{|\mathcal{I}(s_i)| \times |s_j|}$ :

$$C = \text{softmax}(h_{\mathcal{I}(s_i)} \cdot h_{s_j}^\top) \cdot \text{softmax}(h_{s_j} \cdot h_{\mathcal{I}(s_i)}^\top)^\top$$

where softmax is applied on each row. The correspondence matrix  $C$  maps each point in the invariant region  $\mathcal{I}(s_i)$  to the query state  $s_j$ . The matched points in  $s_j$  constitute  $\mathcal{I}(s_j)$  but we only need  $C$  to predict the action pose  $T_j$  of the query state  $s_j$ , as explained in Section III-C2. We detail the training of invariant region matching network in Section III-C4.

2) *Correspondence-based Pose Regression*: The standard practice of 6-DoF pose regression is to obtain the action pose  $T$  from a neural network. However, this approach does not generalize well to new tasks, as shown in Section IV-B. Instead, we propose to analytically compute action pose  $T_j$  of the query state  $s_j$  by solving the optimization problem in Equation 1 with a standard least-squares algorithm [21], as follows,

$$T_j = \arg \min_{T \in \text{SE}(3)} \| T T_i^{-1} P_{\mathcal{I}(s_i)} C - P_{s_j} \|^2, \quad (1)$$

where  $T_i$  is the demonstrated action pose of the support state  $s_i$ ,  $P_{\mathcal{I}(s_i)}$  and  $P_{s_j}$  are the points in  $\mathcal{I}(s_i)$  and  $s_j$ ,  $C$  is the predicted correspondence matrix.  $C$  can be interpreted as an assignment matrix that maps each point in  $\mathcal{I}(s_i)$  to a point in  $s_j$ . Based on Definition III.1, the optimal action pose  $T_j$  is the solution for Equation 1 when the point mapping in correspondence matrix  $C$  produces the minimal overall pairwise distance after applying the transformation  $T_j T_i^{-1}$ . Therefore, the pose regression

problem can be solved by learning to match visual elements with Equation 1. We use the differentiable Procrustes operator from Leopard [27] to solve the least-squares problem. As shown in Section IV-E, our correspondence-based pose regression significantly improves the generalization performance.

It is worth noting that our correspondence-based pose regression shows a resemblance to the point cloud registration (PCR) problem [22]. However, there are three major differences. First, PCR assumes the point clouds are captured from the same scene with different camera viewpoints, while our method aligns point clouds from different scenes and objects. Second, PCR aligns the entire or the majority of two point clouds, while our method finds a sparser matching between only the invariant regions of some given task. Third, PCR finds the best matching, while our method uses the matching as an intermediate step to find the optimal pose of a robot's end-effector in a new scene.

3) *State Routing Network*: We design the state routing network (shown in Figure 4) to select the support frame  $s_i$  in the one-shot demonstration  $\tau$ , given a query scene  $s_j$ . We first extract the scene-level features for the query state  $s_j$  and for each state in  $\tau$  using a PTv2 backbone [10]. Next, we follow the convention of existing work [17] to concatenate the scene-level features with the low-dimensional internal robot state, including joint positions and timesteps. Then, we apply cross-attention over the features of multiple states. Unlike the invariant region matching network, classic attention layers are used instead of graph attention. Finally, we select the state with the strongest attention to the query state  $s_j$  as the support state  $s_i$ .

The three techniques presented above form together the **Invariance Matching One-shot Policy Learning (IMOP)** algorithm. For each new query state  $s_j$ , the support state  $s_i$  is determined from the one-shot demonstration trajectory  $\tau$  by the state routing network. Next, the invariant region matching network predicts the correspondence matrix  $C$  between  $\mathcal{I}(s_i)$  and  $s_j$ . Finally, the action pose  $T_j$  is derived from Equation 1.

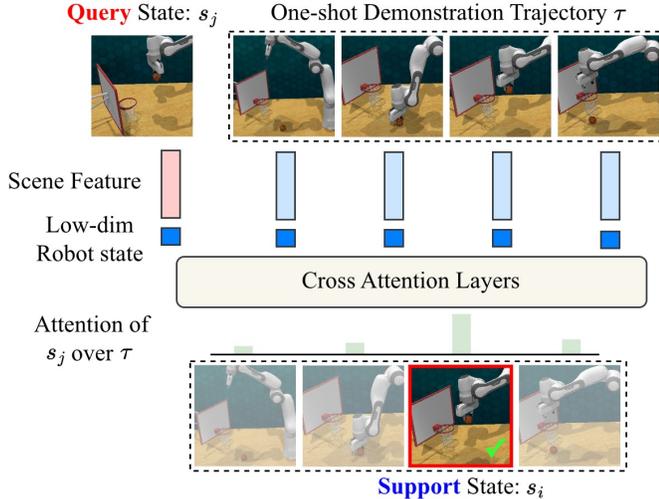


Fig. 4: Overview of the proposed state routing network. Given the query state  $s_j$  and one-shot demonstration trajectory  $\tau$ , we extract the scene-level features for  $s_j$  and each state in  $\tau$ . The scene features are concatenated with the corresponding low-dimensional robot states, and applied with cross-attention layers. The state with the strongest attention to the query state  $s_j$  is selected as the support state  $s_i$ .

We directly use the gripper values  $\lambda_i$  of support state  $s_i$  as the value of  $\lambda_j$  for  $s_j$ . By adapting the demonstration trajectory  $\tau$ , IMOP achieves one-shot generalization on novel tasks.

4) *Training*: To train IMOP, we assume the instance segmentation results are available. We use the segmentation masks provided by RLbench [24]. Note that IMOP does not require any segmentation during inference. Following the convention of RVT [17] and PerAct [38], we represent each state as a macro-step using the keyframe extraction procedure from C2F-ARM [25]. At each training iteration, two trajectories  $\tau_i$  and  $\tau_j$  are sampled if  $|\tau_i| = |\tau_j|$ . The state routing network is trained with focal loss [30] to predict that  $s_i \equiv s_j$  (i.e.,  $s_i$  and  $s_j$  share the same manipulation action) if  $s_i$  and  $s_j$  have the same timestep and vice versa. Let  $\{c_{i,n}\}_{n=1}^N$  denote  $N$  instance segments, where each segment  $c_{i,n} \in s_i$  is a set of 3D points of an object instance in  $s_i$ . To train the invariant region matching network, we estimate the ground-truth  $\mathcal{I}(s_i)$  in Equation 2 as follows,

$$\mathcal{I}(s_i) = \arg \min_{c_{i,n}} \mathbf{E}_{p \in c_{i,n}} \left[ \min_{\substack{q \in c_{j,m} \\ \text{cls}(c_{i,n}) = \\ \text{cls}(c_{j,m})}} \| T_i^{-1}p - T_j^{-1}q \| \right], \quad (2)$$

where  $\text{cls}(c_{i,n})$  denotes the object class of the instance segment  $c_{i,n}$ . The ground-truth  $\mathcal{I}(s_i)$  is the instance segment that has the minimal displacement in the frame of the action pose  $T$ . The ground-truth correspondence between  $s_i$  and  $s_j$  is the set of point mappings  $\{(p, q) \in \mathcal{I}(s_i) \times \mathcal{I}(s_j) | q = \arg \min_{q \in \mathcal{I}(s_j)} \| T_i^{-1}p - T_j^{-1}q \| \}$  such that the point-wise displacement is minimal in the action pose frame.

## IV. EXPERIMENTS

We evaluate **IMOP** on the visual robotic manipulation tasks from RLbench [24], a standard benchmark used in multi-task learning [23]. We first train and evaluate our algorithm on the standard 18 RLbench tasks, and then measure its one-shot generalization ability on 22 novel tasks from 9 categories [19]. Further, we evaluate our method in real-robot experiments through one-shot sim-to-real transfer. Specifically, we aim to answer the following questions:

- Can IMOP achieve high success rates on base tasks?
- Can IMOP generalize to novel tasks given a single demonstration without tuning?
- Can IMOP generalize to objects with large shape variations?
- Can IMOP estimate meaningful invariant regions from demonstrations of novel tasks?
- Can IMOP be trained in simulation and transferred to solve real-robot manipulation tasks?

**Environment.** Following the convention [17], we record  $128 \times 128$  RGB-D images from the front, left/right-shoulders, and wrist cameras. During training, we use 100 recorded trajectories per base task. RLbench contains various task categories, such as pick-and-place, tool use, high-precision operations, screwing, tasks with visual occlusion, and long-term manipulation. We choose 22 novel tasks that have different object setups and task goals from the base ones, according to the task categorization of Hiveformer [19]. For each novel task, only a single successful trajectory is provided, as a one-shot demonstration. Each task is evaluated on 25 independent trials, and we report the average success rate. RLbench generates different object layouts and orientations in different trials of the same task. For the base tasks, we use the same test scenes as the ones in [23]. For novel tasks, we generate 25 test scenes for each task. Figure 2 illustrates the task, and object-level generalizations of IMOP within the RLbench environment. Further details are provided in the supplementary material.

### A. Performance Comparison on Base Tasks

**Baselines.** To evaluate the performance of IMOP on base tasks, we test policies learned from offline demonstrations without the one-shot demonstrations. Specifically, we consider C2F-ARM [25], PerAct [38], RVT [17], and Act3D [16]. C2F-ARM and PerAct learn separate policies for each task. Similar to our IMOP, RVT and Act3D learn a single multi-task policy.

**Results.** Table I compares the success rates of IMOP and other baselines on the base tasks. C2F-ARM and PerAct use 3D voxel representation. RVT renders point clouds to virtual images. Similar to IMOP, Act3D uses raw point clouds as states. Overall, IMOP outperforms all baselines when the best rank and success rate are averaged across all tasks (+4.5% average success rate). IMOP has the highest rank on 10 out of 18 tasks. Remarkably, IMOP achieves a success rate of 52.8% on the “place cups” task while the existing works only reach a rate of 2% to 4%. Moreover, not only it outperforms existing techniques, IMOP also generalizes to novel tasks, while the baselines can only be used on the tasks they were trained on.

TABLE I: Base Task Performance on RL Bench. We report the success rate for each task, and measure the average success rate and rank across all tasks. IMOP has the best overall success rate and rank. The success rate of IMOP is measured with an average of 5 runs. The success rates of baselines are reported from respective papers.

Method	Avg. Success	Avg. Rank	Close Jar	Drag Stick	Insert Peg	Meat off Grill	Open Drawer	Place Cups	Place Wine	Push Buttons
C2F-ARM [25]	20.1	4.9	24	24	4	20	20	0	8	72
PerAct [38]	49.4	3.6	55.2	89.6	5.6	70.4	88	2.4	44.8	92.8
RVT [17]	62.9	2.3	52	<b>99.2</b>	11.2	88	71.2	4	91	<b>100</b>
Act3D [16]	65.1	2.2	<b>92</b>	92	<b>27</b>	<b>94</b>	93	3	80	99
IMOP (Ours)	<b>69.6</b>	<b>1.9</b>	39.2	98.4	12.0	92.0	<b>100.0</b>	<b>52.8</b>	<b>96.0</b>	96.0

	Put in Cupboard	Put in Drawer	Put in Safe	Screw Bulb	Slide Block	Sort Shape	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
C2F-ARM [25]	0	4	12	8	16	8	0	0	0	68
PerAct [38]	28	51.2	84	17.6	74	16.8	26.4	2.4	52	88
RVT [17]	49.6	88.0	91.2	48	81.6	36	28.8	26.4	72	93.6
Act3D [16]	<b>51</b>	90	95	47	<b>93</b>	8	12	9	92	<b>94</b>
IMOP (Ours)	46.4	<b>100.0</b>	<b>96.0</b>	<b>82.0</b>	58.4	<b>37.6</b>	<b>40.0</b>	<b>56.8</b>	<b>100.0</b>	51.2

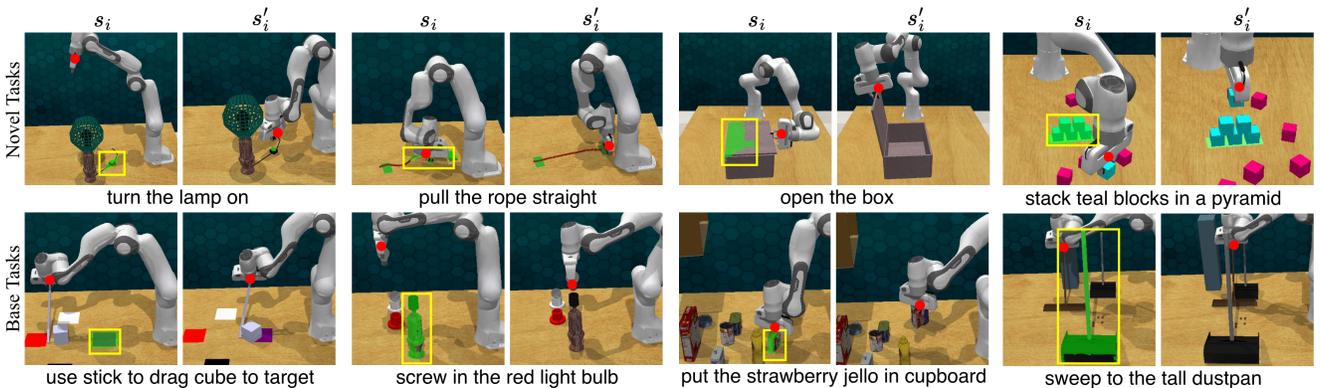


Fig. 5: Visualization of invariant regions  $\mathcal{I}(s_i)$  as estimated by our invariant region matching network on selected base and novel tasks. The invariant regions on  $s_i$  are highlighted with green masks and yellow bounding boxes. The end-effector position is marked with red dot.  $\mathcal{I}(s_i)$  generally covers the target object of the action that is applied by the robot to transition from state  $s_i$  to next state  $s'_i$ . For example, the invariant region of ‘turn the lamp on’ covers the ‘lamp switch’ area. The invariant region of ‘pull the rope straight’ covers one side of the rope and the rope target. It is worth noting that IMOP only takes point cloud as input and has no access to segmentation masks.

### B. Generalization to Novel Tasks

**Baselines.** To assess the one-shot ability of IMOP, and since there is no existing one-shot learning method in the literature that was evaluated in RL Bench, we extend the state-of-the-art technique RVT into two baselines, called RVT+FT and RVT+HDT, that can be used for one-shot learning. Specifically, since fine-tuning is the most common strategy for transferring models to new domains [44], we fine-tune RVT models pretrained on the base tasks with demonstrations of novel tasks and denote this baseline as RVT+FT. HDT [43] is a recent work that achieves one-shot policy generalization by training additional adaptation layers for new tasks while freezing original parameters. We implement the HDT adaptation layers in RVT and denote this baseline model as RVT+HDT. RVT and Act3D are both the most recent methods for learning multi-task robot control policies, and have similar performances, as shown in Table I. We choose RVT over Act3D because the pretrained model of the latter is not yet available. For

both RVT+FT and RVT+HDT, we train the respective network parameters with the one-shot demonstrations of novel tasks until convergence and adopt the sampling and 3D perturbation augmentation strategy from RVT and PerAct.

**Results.** Table II compares the success rates of IMOP, RVT+FT, and RVT+HDT on novel task. IMOP outperforms all baselines with the best rank and success rate averaged across all tasks (+11.5% over RVT+FT). IMOP has the best rank on 15 out of 22 tasks. Despite using the simple fine-tuning strategy, RVT+FT performs better than RVT+HDT on novel tasks. However, we observe the catastrophic forgetting phenomenon [26] in RVT+FT, whose average success rate on base tasks plummets to 4.6% from 62.9% after fine-tuning on novel tasks. RVT+HDT maintains the base task performance by freezing the original parameters but at the cost of having more parameters in the adaptation layers. Parameter sizes are compared in Table III.

Compared to the baselines, IMOP generalizes to novel tasks

TABLE II: Novel Task Performance on RLBench. We report the success rate for each task, and measure the average success rate and rank across all tasks. IMOP has the best overall success rate and rank. Both RVT+FT and RVT+HDT require fine-tuning but IMOP does not. Directly using RVT without fine-tuning results in an average success rate of 0. We do not include this in the table. All success rates are measured with an average of 5 runs.

Method	Avg. Success	Avg. Rank	Ball in Hoop	Rubbish in Bin	Scoop Spatula	Place Hanger	Hit Ball	Block Pyramid	Slide Place	Phone on Base	Open Box	Close Lid
RVT+FT	29.8	2	42.4	50	7.2	4	1.6	0	0	20.4	<b>28</b>	<b>100</b>
RVT+HDT	26.9	2.4	<b>87.6</b>	40.4	0	0	0	0	0	0	24.8	91.2
IMOP (Ours)	<b>41.3</b>	<b>1.5</b>	80.0	<b>94.0</b>	<b>52.0</b>	<b>24.4</b>	<b>32.0</b>	<b>5.2</b>	<b>6.8</b>	<b>22.8</b>	8.0	71.2

	Lid off Pan	Lamp On	Beat Buzz	Remove Cups	Play Jenga	Knife on Board	Straighten Rope	Change Clock	Open Bottle	Open Door	Money out Safe	Close Microwave
RVT+FT	78.8	65.2	19.2	4	40.8	14	0	10.4	34.4	<b>68</b>	22.4	44
RVT+HDT	80.4	<b>72</b>	16.4	0	27.6	8	0	<b>36</b>	23.2	12	12	<b>59.6</b>
IMOP (Ours)	<b>100.0</b>	12.0	<b>20.0</b>	<b>36.0</b>	<b>100.0</b>	<b>20.0</b>	<b>16.8</b>	24.0	<b>81.6</b>	41.2	<b>42.8</b>	21.6

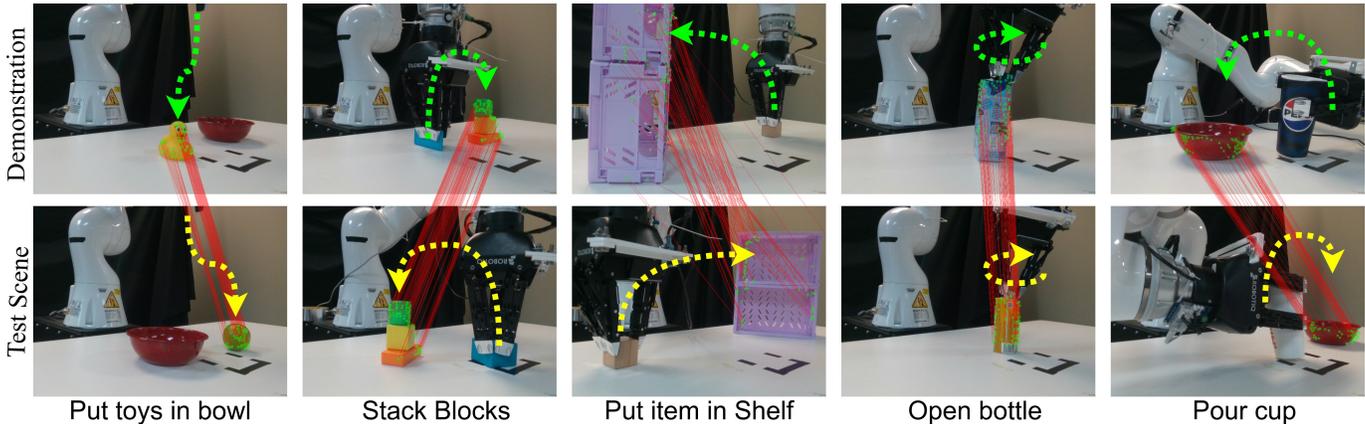


Fig. 6: Real-robot manipulation examples of IMOP via one-shot sim-to-real transfer. The correspondences between the demonstration and test scenes are shown in red lines. The invariant regions closely overlap the next object to interact with. We observe that the one-shot generalization capacity of IMOP, trained with simulation data, enables the sim-to-real transfer using a single real-world demonstration without any re-training.

TABLE III: Network parameters size comparison between IMOP and one-shot baselines.

Method	Success Rate (%)		Params Size
	Base	Novel	
RVT+FT	4.7	29.8	104M
RVT+HDT	62.9	26.9	131M
Ours	69.6	41.3	80.8M

by matching invariant regions between the one-shot demonstration and the states encountered during testing, without further training. Therefore, IMOP does not suffer from catastrophic forgetting or increased size of the parameters. Remarkably, both RVT+FT and RVT+HDT completely fail on certain tasks, *e.g.*, “block pyramid”, “slide place”, and “straighten rope”, while IMOP still delivers a non-trivial performance on these tasks. Both “block pyramid” and “slide place” are categorized as long-horizon tasks and require over 30 and 10 macro-steps, respectively, to finish [19]. The “straighten rope” task requires operations on ropes, which are drastically different from the objects in the base tasks. This further indicates the limitation

of the existing one-shot policy generalization strategies and the advantage of our proposed invariance matching-based policy. Figure 5 visualizes invariant regions on selected base and novel tasks. It can be seen that the invariant regions generally cover the target objects of the robot’s actions. Additional visualizations are provided in the supplementary material.

**Analysis of underperformed tasks.** Despite the overall superiority of IMOP, the performance on specific tasks is still worth discussing. As in Table II, IMOP exhibits a lower success rate on ‘open box’, ‘close microwave’, and ‘close lid’ compared to RVT+FT and RVT+HDT. These tasks involve manipulating flat surfaces, such as lids and box covers, as shown by the invariant region of ‘open box’ visualized in Figure 5. The challenge arises from the potential ambiguity in identifying invariant regions or correspondences on flat surfaces lacking rich geometric features. IMOP also performs worse than the baselines on the ‘lamp on’ task, as you can see in the first two columns of Figure 5. Although the invariant region estimation for this task is good, the area of the switch is small and hence, the correspondence is sparse compared to the whole scene, which leads to unstable solutions of the least-square

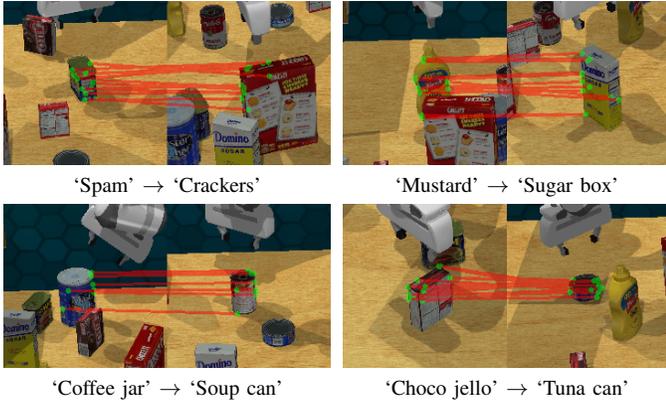


Fig. 7: Visualization of correspondence between different objects for the picking action. ‘Spam’ → ‘Crackers’ denotes using ‘Spam’ as the demonstration to correspond to ‘Crackers’ during testing. Despite large shape variations, IMOP is able to estimate correspondences that effectively transfer actions, as shown by the success rate in Figure 8.

problem in Equation 1. The success rate of IMOP on long-term tasks, *i.e.*, ‘block pyramid’ and ‘slide place’, outperforms the baseline methods but still has a large room for improvement. The IMOP mechanism of routing new states in test scenes to demonstrations is simple yet lacks the failure recovery ability because the demonstrations only contain successful trajectories. A more detailed analysis of failure cases is provided in the supplementary material.

### C. Generalization to Large Shape Variations

**Setup.** Previous experiments evaluate the one-shot generalization ability of IMOP on novel tasks with one demonstration trajectory for each task. Even though there are large variations in object shapes and categories between base and novel tasks, the variations within each novel task are limited. For example, the novel task ‘put the knife on the chopping board’ includes objects such as a knife, knife holder, and chopping board. None of these objects are seen during training and they are placed randomly at each trial and sometimes with color variations, but the object categories and shapes remain unchanged among various trials, as illustrated in Figure 2. We are interested in whether IMOP can generalize to objects with large shape variations. To study this, we evaluate IMOP on objects different from those in demonstrations, which naturally have different shapes. Specifically, we consider a pick-and-place setting where the objective is to pick groceries and place them into a cupboard. The demonstration is given on a certain object, but picking and placing a different object is required during testing. To simplify our study, we assume the access of the segmentation mask of the required object during testing. Therefore, the correspondence between different objects can be obtained by masking out other regions. To avoid memorization, we remove the corresponding state pairs from the training data.

**Result.** Figure 8 shows the success rate of pick-and-place

Demonstration	berry-jello	70	80	10	20	60	70	60	40	30	Success Rate
	choco-jello	60	80	30	20	80	100	70	60	40	
	coffee	40	70	50	20	0	60	40	40	40	
	crackers	20	30	20	30	0	20	30	20	20	
	mustard	60	0	60	20	80	50	0	60	10	
	soup	60	70	60	20	70	90	70	30	20	
	spam	50	50	30	30	30	50	70	30	20	
	sugar	60	50	30	20	60	80	60	60	20	
	tuna	40	40	20	10	0	30	20	20	30	
			berry-jello	choco-jello	coffee	crackers	mustard	soup	spam	sugar	

Fig. 8: Success rate of pick-and-place task using different objects with large shape variations as demonstrations. For example, the entry at the top-right corner, *i.e.*, ‘berry-jello, tuna’, denotes the success rate is 30% at picking and placing ‘tuna can’ while using ‘berry-jello box’ as the demonstration.

with demonstrations on different objects. For example, the entry ‘mustard, coffee’ represents the success rate of picking and placing a coffee jar with a demonstration on a mustard bottle. There are 9 objects in total, and all success rates are measured on an average of 10 trials. It can be seen that IMOP performs competitively when a different object is used as a demonstration, except in some cases of the mustard bottle whose shape is drastically different than others. Using berry-jello and choco-jello as demonstrations delivers the same or even better results on many objects during testing, *e.g.*, soup, spam, sugar, and tuna. These two jello boxes are the smallest objects in this setup.

Figure 7 visualizes the correspondence between the invariant regions of different objects. The corresponding points, visualized as green dots, show a resemblance to the keypoint-based manipulation methods [15, 33]. kPAM [33] and kPAM-SC [15] use object keypoints to compose objectives and constraints in the optimization-based planning for manipulation and demonstrate generalization over large shape variations in a pick-and-place setting. However, these methods require the manual design of keypoints for each object category. Our method IMOP does not require any specification of keypoints or any other form of labeling. The visualized connections are obtained by simply thresholding the correspondence matrix.

### D. Real-Robot Manipulation via One-shot Sim-to-Real Transfer

**Setup.** We evaluate the one-shot sim-to-real performance of IMOP using simulation demonstrations for the base tasks, and trajectories recorded in the real-world as the one-shot demonstrations of novel tasks. We adopt five novel tasks: *put toys in bowl*, *stack blocks*, *put item in shelf*, *open bottle*, and *pour cup*. Contrary to existing work that requires over 10 demonstrations per task for re-training [17], the one-shot

generalization capability of IMOP enables the sim-to-real transfer using only one demonstration and without re-training. We deploy the model on a Kuka LBR iiwa robot. We use RGB-D observations from two RealSense D415 cameras. We use MoveIt [9] for planning paths that move the arm to the gripper poses predicted by IMOP.

TABLE IV: One-shot sim-to-real transfer performance of IMOP for solving real-world manipulation tasks.

Task	Success	
	One-shot RVT	IMOP (Ours)
Put toys in bowl	3/10	8/10
Stack blocks	1/10	4/10
Put item in shelf	2/10	5/10
Open bottle	2/10	4/10
Pour cup	2/10	5/10

**Result.** Table IV shows the success rate of IMOP on the five real-world tasks. The success rate of each task is the average of independent 10 runs with different object layouts. Figure 6 visualizes the correspondence between demonstration and test scenes. For the one-shot RVT, we finetune the pretrained RVT model with a single demonstration for each task. The sim-to-real gap is commonly acknowledged as a challenging problem in deploying control policies trained with simulation data [49]. In our study, we find that the region matching is robust to this domain shift. However, the sigmoid output of the invariant region prediction has much smaller values than those from simulation data. We circumvent this problem by thresholding the sigmoid output with a much smaller value, 0.1. This indicates that IMOP has the one-shot sim-to-real capacity, but still suffers from the sim-to-real gap for the invariant region estimation. This issue can potentially be addressed by fine-tuning with real-robot demonstrations or using more cameras to improve the point cloud quality.

### E. Ablation Studies

We study the component effects of IMOP in Table V. For the state routing network, we find the low-dimensional robot state (i.e., joint positions and timesteps) contributes to the routing accuracy and success rate. This practice is also adopted by many existing work [38, 17, 16].

For the invariant region matching network, we study the attention type, matching strategy, and the degree (K) of the point cloud KNN graphs. Table V shows that KNN graphs of a larger degree provide stronger performance. However, using traditional attention instead of graph attention, which connects all points globally, significantly decreases performance. Note that we only apply traditional attention for region matching and still keep graph attention for invariant region estimation. This suggests local attention and propagating information through KNN graphs are effective inductive biases for 3D region matching. In contrast to the framework in Figure 3, which first estimates  $\mathcal{I}(s_i)$  and then match  $\mathcal{I}(s_i)$  and  $s_j$ , we investigate the strategy of directly matching  $s_i$  and  $s_j$  and predict the correspondence matrix. However, the model fails to predict

TABLE V: Ablation studies on the components of IMOP.

Configurations		Success Rate (%)	
		Base	Novel
<i>State Routing Network</i>			
Using Low-dim robot states	Y	69.6	41.3
	N	63.2	33.0
<i>Invariant Region Matching Network</i>			
Attention Type	Graph Attention	69.6	41.3
	Traditional Attention	52.8	16.4
Matching Strategy	Estimate $\mathcal{I}(s_i)$ , then match $\mathcal{I}(s_i)$ and $s_j$	69.6	41.3
	Directly match $s_i$ and $s_j$	6.7	-
KNN-Graph (K)	4	58.7	32.4
	8	64.5	38.6
	16	69.6	41.3
<i>Regression Method</i>			
Correspondence-based Regression (Equation 1)		69.6	41.3
Regression from matched points		58.2	29.4
Conventional pose regression		65.1	-

actions on all tasks. This indicates the importance of estimating invariant regions.

We compare our correspondence-based pose regression with two alternative solutions. The ‘regression from matched points’ trains extra layers to predict pose from the matched points in correspondence matrix  $C$  instead of deriving the action pose  $T_j$  with Equation 1. The ‘conventional pose regression’ adopts a standard 3D regression head from Act3D [16] while dropping the region matching mechanism completely. Table V shows that conventional regression delivers comparable performance on base tasks but fails to generalize over novel tasks, and deriving action pose with Equation 1 performs better than predicting pose from matched points with an extra regression head.

## V. DISCUSSION AND CONCLUSION

We have shown that one-shot novel task generalization can be achieved by learning to estimate and match key invariant regions in the demonstration and test scene. The target end-effector pose can be transferred by finding correspondences between invariant regions. Some open questions are, however, still worth discussing. IMOP is trained within the RLbench simulation environment, which offers segmentation masks and a validated key frame extraction procedure. For real-world data, a more sophisticated key frame extraction method and the use of open-set object detectors [46, 47, 50] may be necessary. The proposed idea of transferring actions through matching key visual elements is general but the current definition of the invariant region is still closely linked to rigid body transformation. This suggests the potential for extending the formulation of invariant regions based on more general motion descriptors, such as warping or flows [27], to improve the manipulation accuracy on non-rigid objects and tasks with

complex dynamics that cannot be regulated by poses. With the emergence of neural architectures that are scalable to Internet-scale data [6], the possibility of learning invariant regions from unlabeled 2D videos that are widely available has yet to be fully explored. Beyond only leveraging a single demonstration, it is possible to maintain a pool of demonstrations and transfer actions from the most relevant state to improve the manipulation performance and reduce error accumulation under scenes with large variations or requiring failure recovery.

#### REFERENCES

- [1] Paola Ardón, Èric Pairet, Katrin S Lohan, Subramanian Ramamoorthy, and Ronald Petrick. Affordances in robotic tasks—a survey. *arXiv preprint arXiv:2004.07400*, 2020.
- [2] Ondrej Biza, Skye Thompson, Kishore Reddy Pagidi, Abhinav Kumar, Elise van der Pol, Robin Walters, Thomas Kipf, Jan-Willem van de Meent, Lawson LS Wong, and Robert Platt. One-shot imitation learning via interaction warping. *arXiv preprint arXiv:2306.12392*, 2023.
- [3] Abdeslam Boularias, Jens Kober, and Jan Peters. Relative entropy inverse reinforcement learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS*, pages 182–189, 2011.
- [4] Abdeslam Boularias, Oliver Kroemer, and Jan Peters. Learning robot grasping from 3-d images with markov random fields. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS*, pages 1548–1553, 2011.
- [5] Johannes Brandstetter, Rob Hesselink, Elise van der Pol, Erik J Bekkers, and Max Welling. Geometric and physical quantities improve e (3) equivariant message passing. *arXiv preprint arXiv:2110.02905*, 2021.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- [7] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR, 2016.
- [8] Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018.
- [9] David Coleman, Ioan Sucas, Sachin Chitta, and Nikolaus Correll. Reducing the barrier to entry of complex robotic software: a moveit! case study. *arXiv preprint arXiv:1404.3785*, 2014.
- [10] Pointcept Contributors. Pointcept: A codebase for point cloud perception research. <https://github.com/Pointcept/Pointcept>, 2023.
- [11] Sudeep Dasari and Abhinav Gupta. Transformers for one-shot visual imitation. In *Conference on Robot Learning*, pages 2071–2084. PMLR, 2021.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [13] Yan Duan, Marcin Andrychowicz, Bradley Stadie, OpenAI Jonathan Ho, Jonas Schneider, Ilya Sutskever, Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. *Advances in neural information processing systems*, 30, 2017.
- [14] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on robot learning*, pages 357–368. PMLR, 2017.
- [15] Wei Gao and Russ Tedrake. kpam-sc: Generalizable manipulation planning using keypoint affordance and shape completion. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6527–6533. IEEE, 2021.
- [16] Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: 3d feature field transformers for multi-task robotic manipulation. In *Conference on Robot Learning*, pages 3949–3965. PMLR, 2023.
- [17] Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. *arXiv preprint arXiv:2306.14896*, 2023.
- [18] Christian Graf, David B Adrian, Joshua Weil, Miroslav Gabriel, Philipp Schillinger, Markus Spies, Heiko Neumann, and Andras Gabor Kupcsik. Learning dense visual descriptors using image augmentations for robot manipulation tasks. In *Conference on Robot Learning*, pages 871–880. PMLR, 2023.
- [19] Pierre-Louis Guhur, Shizhe Chen, Ricardo Garcia Pinel, Makarand Tapaswi, Ivan Laptev, and Cordelia Schmid. Instruction-driven history-aware policies for robotic manipulations. In *Conference on Robot Learning*, pages 175–187. PMLR, 2023.
- [20] Kai Han, Yunhe Wang, Hanqing Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 45(1):87–110, 2022.
- [21] TS Huang, SD Blostein, and EA Margerum. Least-squares estimation of motion parameters from 3-d point correspondences. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, volume 10, pages 112–115. IEEE Computer Soc. Press Washington DC, 1986.
- [22] Xiaoshui Huang, Guofeng Mei, Jian Zhang, and Rana Abbas. A comprehensive survey on point cloud registration. *arXiv preprint arXiv:2103.02690*, 2021.
- [23] Stephen James and Andrew J Davison. Q-attention: Enabling efficient learning for vision-based robotic manipulation. *IEEE Robotics and Automation Letters*, 7(2): 1612–1619, 2022.
- [24] Stephen James, Zicong Ma, David Rovick Arrojo, and An-

- drew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- [25] Stephen James, Kentaro Wada, Tristan Laidlow, and Andrew J Davison. Coarse-to-fine q-attention: Efficient learning for visual robotic manipulation via discretisation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13739–13748, 2022.
- [26] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [27] Yang Li and Tatsuya Harada. Leopard: Learning partial point cloud matching in rigid and deformable scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5554–5564, 2022.
- [28] Junchi Liang and Abdeslam Boularias. Learning category-level manipulation tasks from point clouds with dynamic graph cnns. In *Proceedings of the 2023 International Conference on Robotics and Automation (ICRA)*, 2023.
- [29] Junchi Liang, Bowen Wen, Kostas E. Bekris, and Abdeslam Boularias. Learning sensorimotor primitives of sequential manipulation tasks from visual demonstrations. In *Proceedings of the 2022 International Conference on Robotics and Automation (ICRA)*, 2022.
- [30] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [31] Clare Lyle, Mark van der Wilk, Marta Kwiatkowska, Yarin Gal, and Benjamin Bloem-Reddy. On the benefits of invariance in neural networks. *arXiv preprint arXiv:2005.00178*, 2020.
- [32] Zhao Mandi, Fangchen Liu, Kimin Lee, and Pieter Abbeel. Towards more generalizable one-shot visual imitation learning. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2434–2444. IEEE, 2022.
- [33] Lucas Manuelli, Wei Gao, Peter Florence, and Russ Tedrake. kpm: Keypoint affordances for category-level robotic manipulation. In *The International Symposium of Robotics Research*, pages 132–157. Springer, 2019.
- [34] Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human videos. *arXiv preprint arXiv:2308.10901*, 2023.
- [35] Kaichun Mo, Leonidas J Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani. Where2act: From pixels to actions for articulated 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6813–6823, 2021.
- [36] Katharina Muelling, Abdeslam Boularias, Betty Mohler, Bernhard Schölkopf, and Jan Peters. Learning strategies in table tennis using inverse reinforcement learning. *Biological Cybernetics*, pages 1–17, 2014. ISSN 0340-1200. doi: 10.1007/s00422-014-0599-1.
- [37] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- [38] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, pages 785–799. PMLR, 2023.
- [39] Elise Van der Pol, Daniel Worrall, Herke van Hoof, Frans Oliehoek, and Max Welling. Mdp homomorphic networks: Group symmetries in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:4199–4210, 2020.
- [40] Xiaoyang Wu, Yixing Lao, Li Jiang, Xihui Liu, and Hengshuang Zhao. Point transformer v2: Grouped vector attention and partition-based pooling. *Advances in Neural Information Processing Systems*, 35:33330–33342, 2022.
- [41] Zhou Xian, Nikolaos Gkanatsios, Theophile Gervet, Tsung-Wei Ke, and Katerina Fragkiadaki. Chaineddiffuser: Unifying trajectory diffusion and keypose prediction for robotic manipulation. In *Conference on Robot Learning*, pages 2323–2339. PMLR, 2023.
- [42] Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang Gan. Prompting decision transformer for few-shot policy generalization. In *international conference on machine learning*, pages 24631–24645. PMLR, 2022.
- [43] Mengdi Xu, Yuchen Lu, Yikang Shen, Shun Zhang, Ding Zhao, and Chuang Gan. Hyper-decision transformer for efficient online policy adaptation. *arXiv preprint arXiv:2304.08487*, 2023.
- [44] Fuchao Yu, Xianchao Xiu, and Yunhui Li. A survey on deep transfer learning and beyond. *Mathematics*, 10(19): 3619, 2022.
- [45] Andy Zeng, Pete Florence, Jonathan Tompson, Stefan Welker, Jonathan Chien, Maria Attarian, Travis Armstrong, Ivan Krasin, Dan Duong, Vikas Sindhwani, et al. Transporter networks: Rearranging the visual world for robotic manipulation. In *Conference on Robot Learning*, pages 726–747. PMLR, 2021.
- [46] Xinyu Zhang and Abdeslam Boularias. Optical flow boosts unsupervised localization and segmentation. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7635–7642. IEEE, 2023.
- [47] Xinyu Zhang, Yuting Wang, and Abdeslam Boularias. Detect every thing with few examples, 2023.
- [48] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16259–16268, 2021.
- [49] Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE symposium series on computational intelligence (SSCI)*, pages 737–744. IEEE,

2020.

- [50] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luwei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16793–16803, 2022.