

ET-SEED: Efficient Trajectory-Level SE(3) Equivariant Diffusion Policy

Chenrui Tie^{1,2*} Yue Chen^{1*} Ruihai Wu^{1*}
Boxuan Dong¹ Zeyi Li^{1,3} Chongkai Gao^{2†} Hao Dong^{1†}
¹Peking University ²National University of Singapore
³University of Chinese Academy of Sciences
gaochongkai@u.nus.edu, hao.dong@pku.edu.cn

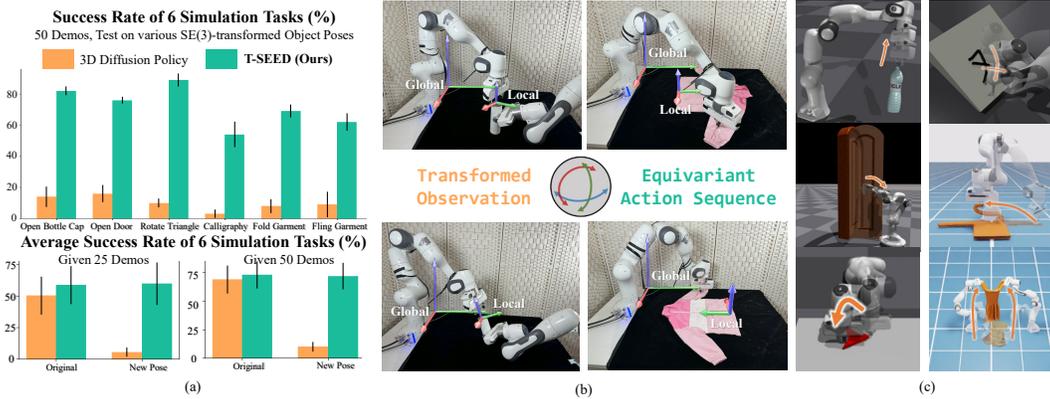


Figure 1: ET-SEED is a visual imitation learning algorithm that marries $SE(3)$ equivariant visual representations with diffusion policies. (a) ET-SEED achieve surprising **efficiency** and **spatial generalization** than baselines. (b) When the input object observation is rotated or translated, the output action sequence change equivariantly. (c) Visualizations of simulation environments.

Abstract: Imitation learning, *e.g.*, diffusion policy, has been proven effective in various robotic manipulation tasks. However, extensive demonstrations are required for policy robustness and generalization. To reduce the demonstration reliance, we leverage spatial symmetry and propose ET-SEED, an efficient trajectory-level $SE(3)$ equivariant diffusion model for generating action sequences in complex robot manipulation tasks. Further, previous equivariant diffusion models require the per-step equivariance in the Markov process, making it difficult to learn policy under such strong constraints. We theoretically extend equivariant Markov kernels and simplify the condition of equivariant diffusion process, thereby significantly improving training efficiency for trajectory-level $SE(3)$ equivariant diffusion policy in an end-to-end manner. We evaluate ET-SEED on representative robotic manipulation tasks, involving rigid body, articulated and deformable object. Experiments demonstrate superior data efficiency and manipulation proficiency of our proposed method, as well as its ability to generalize to unseen configurations with only a few demonstrations. Website: <https://et-seed.github.io/>

Keywords: Imitation Learning, Equivariance, Data Efficiency

1 Introduction

Imitation learning has achieved promising results for acquiring robot manipulation skills [1, 2, 3]. Though, one of the main challenges of imitation learning is that it requires extensive demonstra-

tions to learn a robust manipulation policy [4, 5, 6]. Especially once the spatial pose of the object to be manipulated runs out of the demonstration distribution, the policy performance will easily decrease. Although some works seek to tackle these issues through data augmentation [7] or contrastive learning [8], they usually require task-specific knowledge or extra training, and without theoretical guarantee of spatial generalization ability.

Another promising idea is to leverage symmetry. Symmetry is ubiquitous in the physical world, and many manipulation tasks exhibit a specific type of symmetry known as **SE(3) Equivariance**. $SE(3)$ is a group consisting of 3D rigid transformations. For example, as shown in fig. 1(b), a real robot arm is required to write characters “ICLR” on a paper or fold a garment, when the pose of the paper or the garment changes, the manipulation trajectories of the end-effector should transform equivalently. Employing such symmetries into policy learning can not only improve the data efficiency but also increase the spatial generalization ability. Recent works on 3D manipulation have explored using SE(3) equivariance in the imitation learning process. Most of these works focus on equivariant pose estimation of the target object or end-effector [9, 10, 11]. Trajectory-level imitation learning has achieved state-of-the-art performances on diverse manipulation tasks [3, 12]. By generating a whole manipulation trajectory, this kind of method is capable to tackle more complex manipulation task beyond pick-and-place. For trajectory-level equivariance, Equivariant Diffusion Policy [13] and Equibot [14] propose equivariant diffusion process for robotic manipulation tasks.

However, previous trajectory-level diffusion models for robotic manipulation have two key limitations. First, to maintain equivariance throughout the diffusion process, these models assume that every transition step must preserve equivariance. As we will show in section 3.1, training neural networks with equivariance is more challenging than neural networks with invariance, requiring additional computational resources and leading to slower convergence. This design constrains the model’s efficiency, making it hard for tackling complex long-horizon manipulation tasks. Second, these models define the diffusion process in Euclidean space, which is not a natural definition, and limits the expressiveness. Since the focus is on equivariant diffusion processes within the $SE(3)$ group, it is more natural to define both the diffused variables and the noise as elements of the $SE(3)$ group, which will lead to better convergence and multimodal distributions representation [15].

In this work, we propose **ET-SEED**, a new trajectory-level $SE(3)$ equivariant diffusion model for manipulation tasks. ET-SEED improves the sample efficiency and decreases the training difficulty by restricting the equivariant operations during the diffusion denoising process. We extend the equivariant Markov kernels theory and prove that during the full denoising process, at least only one equivariant transition is required. Then, we integrate the diffusion process on $SE(3)$ manifold [16] and $SE(3)$ transformers [17] to design a new trajectory-level equivariant diffusion model on $SE(3)$ space. In experiment, we evaluate our method on several common and representative manipulation tasks, as shown in fig. 1(c), including rigid body manipulation (rotate triangle, open bottle cap), articulated object manipulation (open door), long-horizon tasks (robot calligraphy), and deformable object manipulation (fold and fling garment). Experiments show our method outperforms SOTA methods in terms of data efficiency, manipulation proficiency and spatial generalization ability (fig. 1(a)). Further, in real-world experiments, with only 20 demonstration trajectories, our method is able to generalize to unseen scenarios.

In summary, our contributions are mainly as followed:

- We propose ET-SEED, an efficient trajectory-level $SE(3)$ equivariant diffusion policy defined on $SE(3)$ manifold, which achieves a proficient and generalizable manipulation policy with only a few demonstrations.
- We extend the theory of equivariant diffusion processes and derive a novel $SE(3)$ equivariant diffusion process, for simplified modeling and inference.
- We extensively evaluate our method on standard robot manipulation tasks in both simulation and real-world settings, demonstrating its data efficiency, manipulation proficiency, and spatial generalization ability, significantly outperforming baseline methods.

2 Related work

2.1 Leveraging Equivariance for Robotic Manipulation

Previous research has demonstrated that leveraging symmetry or equivariance in 3D Euclidean space can improve spatial generalization in a variety of robotic manipulation tasks. Lim et al. [18], Hu et al. [10], Simeonov et al. [19], Xue et al. [20], Gao et al. [11] proposed $SE(3)$ equivariant model for grasp pose prediction. Other works have also leveraged this symmetry in tasks such as part assembly [21, 22], object manipulation on desktop [13], articulated and deformable object manipulation [12] and affordance learning [23]. Most of these studies either focus solely on generating a single 6D pose or fail to guarantee end-to-end equivariance across the entire $SE(3)$ space. In this paper, our proposed method is capable of generating manipulation trajectories while theoretically maintaining end-to-end equivariance over the entire $SE(3)$ group.

2.2 Equivariant Diffusion Models

Diffusion models [24, 25] compose a powerful family of generative models that have proven effective in robotic manipulation tasks [3]. Previous studies [26, 27] have investigated the effectiveness of combining spatial equivariance in the diffusion process to increase data efficiency and improve the spatial generalization ability of the model. GeoDiff [28] gave a theoretical proof of $SE(3)$ equivariant Markov process. Diffusion-EDFs [9] and Orbitgrasp [10] introduced $SE(3)$ equivariant diffusion processes for target grasp pose prediction, but lack the capability to generate entire manipulation trajectories. Wang et al. [13] proposed an equivariant diffusion policy capable of addressing $SO(2)$ equivariant tasks. EquiBot [14] extended equivariant diffusion policies to $SIM(3)$ transformations, with the assumption that every transition step in the diffusion process is equivariant, which demands a high training cost. We further discuss the conditions of $SE(3)$ equivariant diffusion process and prove that not each, but at least one equivariant step is required. Based on this condition, we propose a novel $SE(3)$ equivariant diffusion model achieving better performance than previous works.

2.3 Diffusion on $SE(3)$ manifold

Most diffusion models define the diffusion process on pixel space [25] or 3D Euclidean space [3]. Leach et al. [29] introduces a denoising diffusion model on $SO(3)$ group. $SE(3)$ -Diffusion Fields [15] suggests that in 6-DoF grasp pose generation scenarios, formulating the diffusion process in $SE(3)$ manifold provides better coverage and representation of multimodal distributions, resulting in improved sample efficiency and performance. Jiang et al. [16] proposes a $SE(3)$ diffusion model for robust 6D object pose estimation. In this work, we introduce an equivariant diffusion model on $SE(3)$ manifold for robot manipulation, revealing the superiority of defining equivariant diffusion process on $SE(3)$ over Euclidean space.

3 Method

Problem Formulation. We formulate the problem as an imitation learning setting, aiming to learn a mapping from observation \mathbf{O} to action sequence A , with some demonstrations from an expert policy. In our setting, the observation \mathbf{O} is colored point clouds $\mathbf{P} = \{(x_1, c_1), \dots, (x_N, c_N)\} \in \mathbb{R}^{N \times 6}$. The action is defined directly as the desired 6D pose $\mathbf{H} \in SE(3)$ of the end-effectors. So in our setting, the action sequence means to the trajectory of end-effectors. This experimental setup does not require additional input information, and the action definition is both intuitive and consistent with real robot control, making it applicable to a wide range of robotic manipulation tasks.

In this paper, we propose ET-SEED, a trajectory-level end-to-end $SE(3)$ equivariant diffusion model for robotic manipulation. **ET-SEED can theoretically guarantee the output action is equivariant to any $SE(3)$ transformation applied on the input observation**, while only involving one equivariant denoising step. As shown in fig. 3, given an observation and a noisy action sequence, our

model first implement $K - 1$ invariant denoising steps, and pass the result into the last equivariant denoising step to generate a $SE(3)$ equivariant denoised trajectory.

We will discuss equivariant Markov processes further to explain the correctness and advantages of our proposed diffusion process in section 3.1, with only one denoising step $SE(3)$ equivariant and the rest $SE(3)$ invariant. Then introduce our modified $SE(3)$ invariant and equivariant backbones in section 3.2, and illustrate our $SE(3)$ equivariant diffusion process in section 3.3. Finally we prove end-to-end $SE(3)$ equivariance of our pipeline in section 3.4.

3.1 Equivariant Markov Process

For a Markov process $x^{K:0}$, and any roto-translational transformation $T \in SE(3)$. Geodiff [28] shows that if the initial probabilistic distribution is $SE(3)$ invariant, *i.e.*, $p(x^K) = p(Tx^K)$, and the Markov transitions $p(x^{k-1}|x^k)$ are $SE(3)$ equivariant for any $1 \leq k \leq K$, *i.e.*, $p(x^{k-1}|x^k) = p(Tx^{k-1}|Tx^k)$, then the density of x_0 satisfies $p(x_0) = p(Tx_0)$. Equibot [14] adapts the theory and makes it more consistent with the robotics setting. They involve an additional condition c (can be seen as an observation) and show that if the initial distribution $p(x^K|c)$ and transitions are all equivariant, *i.e.*, $p(x^K|c) = p(Tx^K|Tc)$, $p(x^{k-1}|x^k, c) = p(Tx^{k-1}|Tx^k, Tc)$ then the marginal distribution satisfies $p(x^0|c) = p(Tx^0|Tc)$.

In this paper, we discover that the condition of getting an equivariant marginal distribution $p(x^0|c)$ can be weaker. Formally, we first define three Markov transitions with different properties.

$$\begin{aligned} p_1(x^{k-1}|x^k, c) &= p_1(x^{k-1}|x^k, Tc) \\ p_2(x^{k-1}|x^k, c) &= p_2(Tx^{k-1}|x^k, Tc) \\ p_3(x^{k-1}|x^k, c) &= p_3(Tx^{k-1}|Tx^k, Tc) \end{aligned} \quad (1)$$

Then we derive the marginal distribution using the three types of Markov transitions. We have the following statement. See appendix C for the detailed proof.

Proposition 1 *For a Markov process $x^{K:0}$, if the initial distribution $p(x^K|c) = p(x^K|Tc)$, first $K - n + 1$ transitions follow the property of p_1 , the middle 1 transition follows p_2 , and the last $n - 2$ transitions follow p_3 , then the final marginal distribution satisfies $p(x^0|c) = p(Tx^0|Tc)$.*

Previous works [13, 14] make all transitions p_3 -like, which is a special case of proposition 1. In practice, we observe that training neural networks to approximate the properties of p_2 and p_3 is much more challenging compared to p_1 , both in terms of final performance and training cost. When the condition c is transformed by a $SE(3)$ element, the distributions in p_2 and p_3 change equivalently, while the distribution in p_1 remains unchanged. Learning to output an equivariant feature is clearly more challenging for neural networks than producing an invariant feature. Additionally, in most of implementations of equivariant networks, building and training a model whose output is $SE(3)$ equivariant to the input takes up more computing resources than a $SE(3)$ invariant version. We design experiments to validate these facts, with details in appendix E.

In ET-SEED, we set the parameter $n = 2$, meaning there are $K - 1$ p_1 -like transitions (referred to as “ $SE(3)$ Invariant Denoising Steps”) and one p_2 -like transition (referred to as the “ $SE(3)$ Equivariant Denoising Step”). This key design choice significantly reduces the training complexity, thereby enhancing the overall performance of our method.

3.2 $SE(3)$ equivariant backbone

In order to generate whole manipulation trajectories, it’s necessary that the network has the ability to output a translation vector at anywhere in the 3D space (even beyond the convex hull of the object), which can not be achieved by directly using existing equivariant backbones [17, 30, 31]. In this paper, based on $SE(3)$ Transformer [17], we propose $SE(3)$ equivariant and invariant backbones suitable for predicting $SE(3)$ action sequences, which theoretically satisfy definition 1 and 2

. The implementation details can be found at appendix F. The input of backbone is a set of points coordinates $\mathcal{X} \in \mathbb{R}^{N \times 3}$, with some type-0 features D_0^{in} and type-1 features D_1^{in} attached on each point. Type-0 vectors are invariant under roto-translation transformations and type-1 vectors rotate and translate according to $SE(3)$ transformation of point cloud. The output is an element of $SE(3)$, represented as a 4×4 matrix. For the $SE(3)$ equivariant model (denoted as \mathcal{E}_{equiv}), we have

$$\forall T \in SE(3) : T\mathcal{E}_{equiv}(\mathcal{X}; D_0^{in}, D_1^{in}) = \mathcal{E}_{equiv}(T\mathcal{X}; D_0^{in}, TD_1^{in}) \quad (2)$$

And for the $SE(3)$ invariant model (denoted as \mathcal{E}_{inv}), we have

$$\forall T \in SE(3) : \mathcal{E}_{equiv}(\mathcal{X}; D_0^{in}, D_1^{in}) = \mathcal{E}_{inv}(T\mathcal{X}; D_0^{in}, TD_1^{in}) \quad (3)$$

Algorithm 1 Training phase

```

repeat
  Sample  $A^0, O \sim p_{data}$ 
  Sample  $k \sim \text{Uniform}(\{1, \dots, K\})$ 
  for  $\mathbf{H}_i \in A^0$  do
    Sample  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
     $\mathbf{H}_i^k = \text{Exp}(\gamma\sqrt{1-\bar{\alpha}_t}\varepsilon) \mathcal{F}(\sqrt{\bar{\alpha}_t}; \mathbf{H}^0, \mathbb{H})$ 
    Predict  $\hat{\mathbf{H}}_i^{k \rightarrow 0} = s_\theta(\mathbf{O}, \mathbf{H}_i^k; \mathbf{k}, \mathbf{i})$ 
  end for
  Optimize loss  $\mathcal{L} = \text{loss}(\hat{\mathbf{H}}_i^{k \rightarrow 0}, \mathbf{H}_i(\mathbf{H}_i^k)^{-1})$ 
until converged

```

Algorithm 2 Inference phase

```

for  $\mathbf{H}_i^K \in A^K$  do
  for  $k = K, \dots, 2$  do
    Predict  $\hat{\mathbf{H}}_i^{k \rightarrow 0} = \mathcal{E}_{inv}(O, \mathbf{H}_i^k; k, i)$ 
    Update  $\mathbf{H}_i^{k-1} = \text{Exp}(\lambda_0 \text{Log}(\hat{\mathbf{H}}_i^{k \rightarrow 0} \mathbf{H}_i^k) + \lambda_1 \text{Log}(\mathbf{H}_i^k))$ 
  end for
  Predict  $\hat{\mathbf{H}}_i^{1 \rightarrow 0} = \mathcal{E}_{equiv}(O, \mathbf{H}_i^k; 1, i)$ 
  Assign  $\mathbf{H}_i^0 = \hat{\mathbf{H}}_i^{1 \rightarrow 0} \mathbf{H}_i^k$ 
end for
Return:  $A^0 = \bigcup_i \mathbf{H}_i^0$ 

```

3.3 $SE(3)$ Diffusion Models

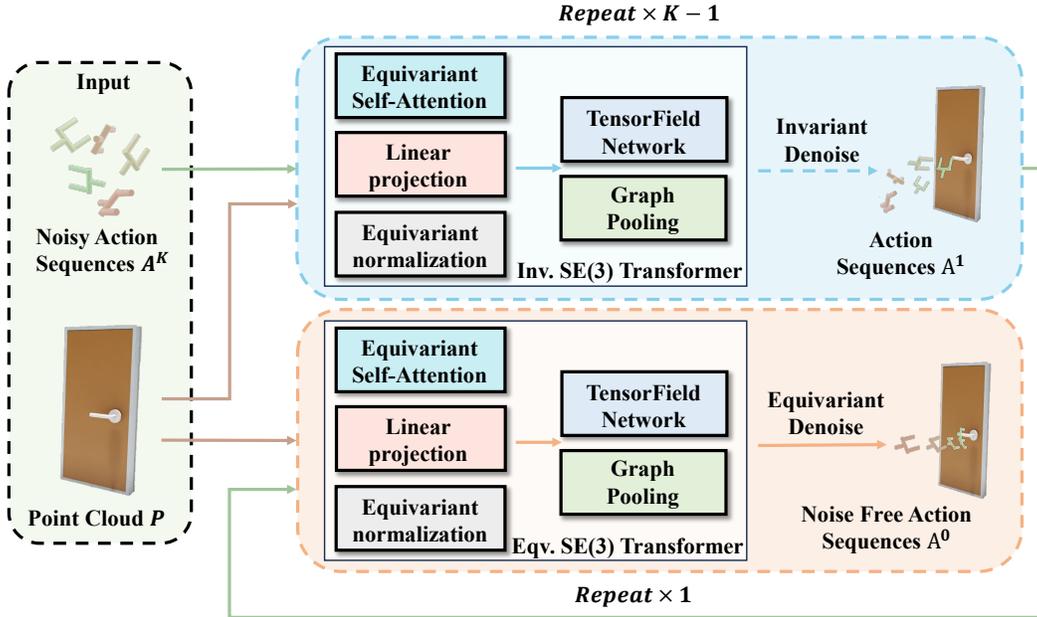


Figure 2: Overview of our pipeline. A colored point cloud and a random sampled action sequence are first passed through $K - 1$ $SE(3)$ invariant denoising steps and then a $SE(3)$ equivariant denoising step to generate a noise free action sequence.

Inspired by standard diffusion model, ET-SEED progressively disturbs the action \mathbf{H}^0 into a noisy action \mathbf{H}^K . As standard diffusion process assume the final noisy variable x_T follows the standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$, we assume the noisy action \mathbf{H}^k follow a Gaussian distribution on $SE(3)$, centered at the identity transformation \mathbb{H} . So we use an interpolation-based $SE(3)$ diffusion formula, which represent the $\mathbf{H}^k \sim q(\mathbf{H}^k | \mathbf{H}^0)$ at noise step $k (1 \leq k \leq K)$ as

$$\mathbf{H}^k = \underbrace{\text{Exp}(\gamma\sqrt{1-\bar{\alpha}_t}\varepsilon)}_{\text{Perturbation}} \underbrace{\mathcal{F}(\sqrt{\bar{\alpha}_t}; \mathbf{H}^0, \mathbb{H})}_{\text{Interpolation}}, \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (4)$$

The interpolation function $\mathcal{F}(\sqrt{\bar{\alpha}_t}; \mathbf{H}^0, \mathbb{H})$ is an intermediate transformation between the origin action \mathbf{H}^0 and the identity transformation \mathbb{H} . By adding a perturbation noise $\text{Exp}(\gamma\sqrt{1-\bar{\alpha}_t}\varepsilon)$ on the intermediate transformation, we get a diffused action \mathbf{H}^k . The training process of ET-SEED is shown in algorithm 1. More explanation about 4 can be found in Jiang et al. [16] or appendix G.

This formulation is an analogy of DDPM [25], which represent the noisy image as

$$x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \bar{\varepsilon}, \bar{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (5)$$

We can treat the first term of 5 as interpolation between x_0 and 0, second term as external noise.

The goal of $SE(3)$ reverse process is to train a denoising network, gradually refine the noisy action to the optimal ones. In this paper, we propose a novel $SE(3)$ equivariant denoising process, which theoretically keeps the equivariance property while reducing the training difficulty. As illustrated in fig. 2, the input of reverse process is an observation O , a noisy action sequence $A^K = [\mathbf{H}_0^K, \mathbf{H}_1^K, \dots, \mathbf{H}_{T_p}^K]$, where T_p is the action prediction horizon, and each \mathbf{H}_i^K is drawn from a $SE(3)$ Gaussian distribution centered at identity transformation \mathbb{H} . The denoising process forms a reverse Markov chain $A^K \rightarrow A^{K-1} \rightarrow \dots \rightarrow A^0$. In our setting, we design a model to predict a probability distribution of the action sequence A^{k-1} given current observation O and the action sequence A^k of the last denoising step. So the input of our denoising network s_θ consists of observation O , noisy action \mathbf{H}_i^k , and scalar condition k, i , outputs the predicted relative transformation between \mathbf{H}_i^k and ground truth. Formally, we have

$$\hat{\mathbf{H}}_i^{k \rightarrow 0} = s_\theta(O, \mathbf{H}_i^k; \mathbf{k}, \mathbf{i}) \quad (6)$$

To ensure the overall $SE(3)$ equivariance of our pipeline, we propose a novel design of denoising network s_θ . It consists of one $SE(3)$ invariant backbone \mathcal{E}_{inv} and one $SE(3)$ equivariant backbone \mathcal{E}_{equiv} . We use \mathcal{E}_{inv} to denoise in the first $K-1$ iterations and use \mathcal{E}_{equiv} to denoise in the last iteration. Formally, our denoising network s_θ is

$$s_\theta(O, \mathbf{H}_i^k; \mathbf{k}, \mathbf{i}) = \begin{cases} \mathcal{E}_{inv}(O, \mathbf{H}_i^k; \mathbf{k}, \mathbf{i}), & \mathbf{k} > \mathbf{1} \\ \mathcal{E}_{equiv}(O, \mathbf{H}_i^k; \mathbf{k}, \mathbf{i}), & \mathbf{k} = \mathbf{1} \end{cases} \quad (7)$$

As illustrated in algorithm 2, in the first $K-1$ denoising iteration, we use $SE(3)$ invariant backbone \mathcal{E}_{inv} to predict noise, and implement a denoise step by

$$\mathbf{H}_i^{k-1} = \text{Exp}(\lambda_0 \text{Log}(\hat{\mathbf{H}}_i^{k \rightarrow 0} \mathbf{H}_i^k) + \lambda_1 \text{Log}(\mathbf{H}_i^k)) \quad (8)$$

This formulation, by minimizing the KL divergence of the posterior distribution and prior distribution of \mathbf{H}_i^{k-1} , is able to infer a more reliable distribution for \mathbf{H}_i^{k-1} [16].

In the last denoising iteration, we use a $SE(3)$ equivariant backbone \mathcal{E}_{equiv} to predict noise and directly apply the predicted transformation on current action.

$$\mathbf{H}_i^0 = \hat{\mathbf{H}}_i^{1 \rightarrow 0} \mathbf{H}_i^k \quad (9)$$

3.4 Proof of end-to-end Equivariance

In this paper, we derive equivariant Markov process from a weaker condition, and build up a practically feasible framework. In this part, we will prove that our proposed denoising process is theoretically $SE(3)$ equivariant. When the input observation O is transformed by any $SE(3)$ element T , the output denoised action sequence A^0 will be equivariantly transformed. Denote the denoising process as $A^0 = ETSEED(A^K; O)$, we have the following statement

Proposition 2 *For a Markov process $A^{K:0}$, if A^K is sampled from Gaussian distribution, $A^0 = ETSEED(A^K; O)$. Then, $\forall T \in SE(3) : TA^0 = ETSEED(TO, A^K)$.*

A detailed proof is attached in appendix D. Here we briefly introduce the intuition. In the first $K-1$ denoising steps, as we use $SE(3)$ invariant backbone to predict noise, the predicted noise

$\hat{H}_i^{k \rightarrow 0}$ is $SE(3)$ invariant, thus the updated H_i^{k-1} is also invariant. After the first $K - 1$ steps, we get an invariant H_i^1 . Because of the equivariance property of our last noise prediction backbone, the predicted noise of last step $\hat{H}_i^{1 \rightarrow 0}$ is $SE(3)$ equivariant. Carried into 9, we get a $SE(3)$ equivariant H_i^0 for all $0 \leq i \leq T_p$. It means the final denoised action sequence A^0 is $SE(3)$ equivariant to the input observation O .

4 Experiments

We systematically evaluate ET-SEED through both simulation and real-world experiments, aiming to address the following research questions: (1) Does our method demonstrate superior spatial generalization compared to existing imitation learning approaches? (2) Can our method achieve comparable performance with fewer demonstrations? (3) Is our method applicable to real-world robotic manipulation tasks?

4.1 Simulation Experiments

Tasks. We design six representative robot manipulation tasks: *Open Bottle Cap*, *Open Door*, *Rotate Triangle*, *Calligraphy*, *Cloth Folding*, and *Cloth Fling*. These tasks encompass manipulation of rigid bodies, articulated bodies, and deformable objects, as well as dual-arm collaboration, long-horizon tasks, and complex manipulation scenarios. A brief overview is illustrated in fig. 1. For each task, we set up multiple cameras to capture full point clouds of the objects to be manipulated. We assume each robot manipulator operates within a complete 6DoF $SE(3)$ action space. Further experiments details and discussions of their equivariant properties can be found in appendix H.

Baselines. We compare our method against the following baselines:

- **3D Diffusion Policy (DP3)** [32]: A diffusion-based 3D visuomotor policy.
- **3D Diffusion Policy with Data Augmentations (DP3+Aug)**: Same architecture as DP3, with $SE(3)$ data augmentation added.
- **EquiBot** [14]: A baseline combines SIM(3)-equivariant neural network architectures with diffusion policy.

DP3 and DP3+Aug are used to compare ET-SEED with baseline methods that utilize data augmentation to achieve spatial generalization, while EquiBot allows for a comparison between different architectures of equivariant diffusion process.

Augmentations. The DP3+Aug baseline utilizes augmentations during training. In all environments, training data is augmented by (1) rotating the observation along all three axes by random angles between 0° and 90° , and (2) applying a random Gaussian offset to the observation. The standard deviation of the Gaussian noise is set to 10% of the workspace size.

Evaluation. Following the setup of Gao et al. [11], we collect demonstrations and train our policy under the **Training setting (T)**, subsequently testing the trained policy on both T and **New Poses (NP)**, where target object poses undergo random $SE(3)$ transformations. We evaluate all methods using two metrics, based on 20 evaluation rollouts, averaged over 5 random seeds. Since we generate complete manipulation trajectories, the final success rate alone is inadequate for fully assessing the trajectory’s quality. We calculate the geodesic distance between each step of the predicted trajectory and the ground truth trajectory, providing a more comprehensive reflection of the trajectory’s overall quality. The geodesic distance between each step of the predicted trajectory and the ground truth trajectory, we can obtain a more accurate reflection of the trajectory’s overall quality. The definition of geodesic distance between $T, \hat{T} \in SE(3)$ is

$$\mathcal{D}_{geo}(T, \hat{T}) = \sqrt{\left\| \text{Log}(\mathbf{R}^\top \hat{\mathbf{R}}) \right\|^2 + \|\hat{\mathbf{t}} - \mathbf{t}\|^2}, \quad (10)$$

where R and t are the rotation and translation parts of T . We report \mathcal{D}_{geo} in the same manner as success rates.

Results. Table 2 and 3 provide a quantitative comparison between our method and the baseline. Both DP3 and its augmented variant demonstrate strong performance in the training setting (T), but they exhibit a significant drop in performance when faced with New Poses (NP) scenarios. This highlights that merely incorporating data augmentation is insufficient for the model to generalize effectively to unseen poses. Instead, leveraging equivariance proves essential for enhancing spatial generalization.

While EquiBot achieves commendable results in both success rate and \mathcal{D}_{geo} , it struggles with more complex, long-horizon tasks such as Calligraphy and Fold Garment. Also, when less demonstrations are given, the performance is not satisfactory. These challenges stem from the inherent complexity of its diffusion process design, where maintaining equivariance in each Markov transition adds substantial difficulty to the learning task.

In contrast, ET-SEED consistently outperforms across all six tasks, with minimal performance drop when facing unseen object poses. This advantage is especially pronounced when using a limited number of demonstrations, showcasing T-SEED’s superior data efficiency, manipulation proficiency, and spatial generalization ability.

4.2 Real-Robot Experiment

Setup. We test the performance of our model on four tasks on real scenarios. All the tasks are visualized in Figure 5. We use Segment Anything Model 2 (SAM2) [33] to segment the object from the scene and project the segmented image with depth to point cloud. Please refer to appendix I and our website for more details and videos of real-world manipulations.

Expert demonstrations are collected by human tele-operation. The Franka arm and the gripper are teleoperated by the keyboard. Since our tasks contain more than one stage and include two robots and various objects, making the process of demonstration collection very time-consuming, we only provide 20 demonstrations for each task.

In test setting, We place the object at 10 different positions with different poses that are unseen in the training data. Each position is evaluated with one trial.

Results. Results for our real robot tasks are given in Table 4. Consistent with our simulation findings, in real world experiments, ET-SEED performs better than baselines in all the four tasks, given only 20 demonstrations. The evaluation shows the effectiveness and spatial generalization ability of our method.

5 Conclusion

In this paper, we propose ET-SEED, an efficient trajectory-level $SE(3)$ equivariant diffusion policy. By leveraging $SE(3)$ symmetries, our method enhances both data efficiency and spatial generalization while reducing the training complexity typically encountered in diffusion-based methods. Through theoretical extensions of equivariant Markov kernels, we demonstrated that the $SE(3)$ equivariant diffusion process can be achieved given a weaker condition, significantly simplifying the learning task. Experimental results on diverse robotic manipulation tasks show that ET-SEED performs better than SOTA methods. Real-world experiments further validate the generalization ability of our model to unseen object poses with only 20 demonstrations. Our work presents a novel approach for data efficient and generalizable imitation learning, paving the way for more capable and adaptive robots in real-world applications.

References

- [1] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation with object proposal priors. In *Conference on Robot Learning*, pages 1199–1210. PMLR, 2023.
- [2] Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.
- [3] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.
- [4] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [5] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [6] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
- [7] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, J. Peralta, B. Ichter, et al. Scaling robot learning with semantically imagined experience. *arXiv preprint arXiv:2302.11550*, 2023.
- [8] T. Ma, J. Zhou, Z. Wang, R. Qiu, and J. Liang. Contrastive imitation learning for language-guided multi-task robotic manipulation. *arXiv preprint arXiv:2406.09738*, 2024.
- [9] H. Ryu, J. Kim, H. An, J. Chang, J. Seo, T. Kim, Y. Kim, C. Hwang, J. Choi, and R. Horowitz. Diffusion-edfs: Bi-equivariant denoising generative modeling on $se(3)$ for visual robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18007–18018, 2024.
- [10] B. Hu, X. Zhu, D. Wang, Z. Dong, H. Huang, C. Wang, R. Walters, and R. Platt. Orbitgrasp: $se(3)$ -equivariant grasp learning. *arXiv preprint arXiv:2407.03531*, 2024.
- [11] C. Gao, Z. Xue, S. Deng, T. Liang, S. Yang, L. Shao, and H. Xu. Riemann: Near real-time $se(3)$ -equivariant robot manipulation without point cloud segmentation, 2024. URL <https://arxiv.org/abs/2403.19460>.
- [12] J. Yang, C. Deng, J. Wu, R. Antonova, L. Guibas, and J. Bohg. Equivact: Sim(3)-equivariant visuomotor policies beyond rigid object manipulation. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9249–9255. IEEE, 2024.
- [13] D. Wang, S. Hart, D. Surovik, T. Kelestemur, H. Huang, H. Zhao, M. Yeatman, J. Wang, R. Walters, and R. Platt. Equivariant diffusion policy, 2024.
- [14] J. Yang, Z. ang Cao, C. Deng, R. Antonova, S. Song, and J. Bohg. Equibot: Sim(3)-equivariant diffusion policy for generalizable and data efficient learning. *arXiv preprint arXiv:2407.01479*, 2024.
- [15] J. Urain, N. Funk, J. Peters, and G. Chalvatzaki. $Se(3)$ -diffusionfields: Learning smooth cost functions for joint grasp and motion optimization through diffusion. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5923–5930. IEEE, 2023.
- [16] H. Jiang, M. Salzmann, Z. Dang, J. Xie, and J. Yang. $Se(3)$ diffusion model-based point cloud registration for robust 6d object pose estimation. *Advances in Neural Information Processing Systems*, 36, 2024.

- [17] F. Fuchs, D. Worrall, V. Fischer, and M. Welling. Se (3)-transformers: 3d roto-translation equivariant attention networks. *Advances in neural information processing systems*, 33:1970–1981, 2020.
- [18] B. Lim, J. Kim, J. Kim, Y. Lee, and F. C. Park. Equigraspflow: Se (3)-equivariant 6-dof grasp pose generative flows. In *8th Annual Conference on Robot Learning*, 2024.
- [19] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitzmann. Neural descriptor fields: Se (3)-equivariant object representations for manipulation. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6394–6400. IEEE, 2022.
- [20] Z. Xue, Z. Yuan, J. Wang, X. Wang, Y. Gao, and H. Xu. Useek: Unsupervised se (3)-equivariant 3d keypoints for generalizable manipulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1715–1722. IEEE, 2023.
- [21] R. Wu, C. Tie, Y. Du, Y. Zhao, and H. Dong. Leveraging se (3) equivariance for learning 3d geometric shape assembly. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14311–14320, 2023.
- [22] G. Scarpellini, S. Fiorini, F. Giuliani, P. Morerio, and A. Del Bue. Diffassemble: A unified graph-diffusion model for 2d and 3d reassembly. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024.
- [23] Y. Chen, C. Tie, R. Wu, and H. Dong. Eqvafford: Se (3) equivariance for point-level affordance learning. *arXiv preprint arXiv:2408.01953*, 2024.
- [24] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [25] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [26] J. Guan, W. W. Qian, X. Peng, Y. Su, J. Peng, and J. Ma. 3d equivariant diffusion for target-aware molecule generation and affinity prediction. *arXiv preprint arXiv:2303.03543*, 2023.
- [27] A. Schneuing, Y. Du, C. Harris, A. Jamasb, I. Igashov, W. Du, T. Blundell, P. Lió, C. Gomes, M. Welling, et al. Structure-based drug design with equivariant diffusion models. *arXiv preprint arXiv:2210.13695*, 2022.
- [28] M. Xu, L. Yu, Y. Song, C. Shi, S. Ermon, and J. Tang. Geodiff: A geometric diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.
- [29] A. Leach, S. M. Schmon, M. T. Degiacomi, and C. G. Willcocks. Denoising diffusion probabilistic models on SO(3) for rotational alignment. In *ICLR 2022 Workshop on Geometrical and Topological Representation Learning*, 2022. URL <https://openreview.net/forum?id=BY88eBbkpe5>.
- [30] C. Deng, O. Litany, Y. Duan, A. Poulernard, A. Tagliasacchi, and L. J. Guibas. Vector neurons: A general framework for so (3)-equivariant networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12200–12209, 2021.
- [31] Y.-L. Liao and T. Smidt. Equiformer: Equivariant graph attention transformer for 3d atomistic graphs. *arXiv preprint arXiv:2206.11990*, 2022.
- [32] Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. *arXiv preprint arXiv:2403.03954*, 2024.

- [33] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. URL <https://arxiv.org/abs/2408.00714>.
- [34] H. Xue, Y. Li, W. Xu, H. Li, D. Zheng, and C. Lu. Unifolding: Towards sample-efficient, scalable, and generalizable robotic garment folding, 2023.
- [35] A. Canberk, C. Chi, H. Ha, B. Burchfiel, E. Cousineau, S. Feng, and S. Song. Cloth funnels: Canonicalized-alignment for multi-purpose garment manipulation, 2022.
- [36] R. Wu, H. Lu, Y. Wang, Y. Wang, and H. Dong. Unigarmentmanip: A unified framework for category-level garment manipulation via dense visual correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16340–16350, June 2024.
- [37] H. Lu, R. Wu, Y. Li, S. Li, Z. Zhu, C. Ning, Y. Shen, L. Luo, Y. Chen, and H. Dong. Garment-lab: A unified simulation and benchmark for garment manipulation. In *Advances in Neural Information Processing Systems*, 2024.

A Preliminary Background

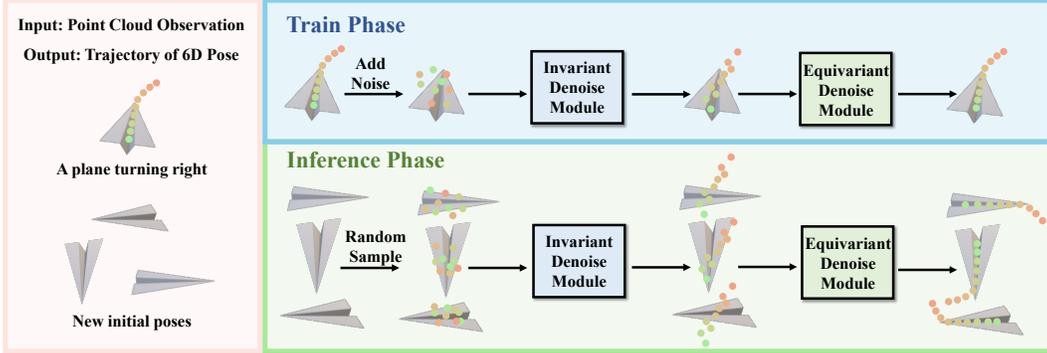


Figure 3: Left: Given a point cloud observation, the model outputs a trajectory of 6D pose. Right: For a noisy trajectory, an $SE(3)$ invariant denoise module generates $SE(3)$ invariant trajectories, followed by an $SE(3)$ equivariant module to produce noise free equivariant trajectories on the $SE(3)$ manifold. The model can generate corresponding trajectories even for unseen poses.

$SE(3)$ Group and its Lie Algebra. $SE(3)$ (Special Euclidian Group) is a group consisting of 3D rigid transformations. Each $SE(3)$ transformation can be represented as a 4×4 matrix (denoted as T), indicating linear transformation on homogeneous 4-vectors. Formally, $T = \begin{bmatrix} R & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$, where $R \in \mathbb{R}^{3 \times 3}$ is a rotation matrix, $\mathbf{t} \in \mathbb{R}^3$ is a translation vector. The Lie algebra $\mathfrak{se}(3)$ is a linear 6D vector space corresponding to the tangent space of $SE(3)$. Each element of $\mathfrak{se}(3)$ is a 6D vector $\delta \in \mathbb{R}^6$. The mutual mapping between $SE(3)$ and $\mathfrak{se}(3)$ is achieved by the logarithm map $\text{Log} : SE(3) \rightarrow \mathbb{R}^6$ and the exponential map $\text{Exp} : \mathbb{R}^6 \rightarrow SE(3)$. More information about $SE(3)$ and $\mathfrak{se}(3)$ can be found in appendix B.

$SE(3)$ Equivariant Function. Generally, we call a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ that maps elements from input space \mathcal{X} to output space \mathcal{Y} is equivariant to a group G if there are group representations of G on \mathcal{X} and \mathcal{Y} respectively denoted by $\rho^{\mathcal{X}}$ and $\rho^{\mathcal{Y}}$ such that $\forall_{g \in G} : \rho^{\mathcal{Y}}(g) \circ f = f \circ \rho^{\mathcal{X}}(g)$

In other words, the function f commutes with representations of the group G . In this paper, we focus on $SE(3)$ group, and represent elements of $SE(3)$ as 4×4 matrixes. As special cases of general equivariance, we can define $SE(3)$ equivariant and invariant functions as:

Definition 1 $SE(3)$ Equivariant Function.

A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is called $SE(3)$ equivariant if

$$\forall_{T \in SE(3)} : T \circ f = f \circ T \quad (11)$$

Definition 2 $SE(3)$ Invariant Function.

A function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is called $SE(3)$ invariant if

$$\forall_{T \in SE(3)} : f = f \circ T \quad (12)$$

$SE(3)$ Equivariant Trajectory. In many robotic manipulation tasks, the trajectories of the manipulator show a certain symmetry. If the representation of a trajectory under certain coordinate frame is $SE(3)$ invariant, we call the trajectory as $SE(3)$ Equivariant. Formally, it can be defined as

Definition 3 $SE(3)$ Equivariant Trajectory.

A trajectory $\{s_i\}_{i=1}^n$ is called $SE(3)$ equivariant if exists a coordinate frame \mathcal{A} , such that for any transformation $T \in SE(3)$ applied on both the trajectory and the coordinate frame (denoted as $\{s'_i\}_{i=1}^n = T\{s_i\}_{i=1}^n$ and $\mathcal{A}' = T\mathcal{A}$), the representation of $\{s'_i\}_{i=1}^n$ by the basis of \mathcal{A}' is same as the representation of $\{s_i\}_{i=1}^n$ by the basis of \mathcal{A} .

This property means the trajectory is “attached” on a certain frame, and when the frame transforms, the trajectory transforms accordingly. To demonstrate the universality of such symmetrical trajectories, we use a general example in fig. 3 . For planes at different position and orientation, the representations of “turning right” trajectories appear identical in each plane’s body coordinate frame. In our experiments, we select 6 representative manipulation tasks showing this symmetry. Further discussion can be found in appendix H .

B $SE(3)$ group and $se(3)$ algebra

B.1 Basic Terms

An element in $SE(3)$ can be represented as a 4×4 matrix

$$\begin{aligned} \mathbf{R} &\in SO(3), t \in \mathbb{R}^3 \\ C &= \begin{pmatrix} \mathbf{R} & t \\ 0 & 1 \end{pmatrix} \end{aligned} \quad (13)$$

This representation means we can compute the composition and inversion of elements in $SE(3)$ by matrix multiplication and inversion.

An element $\delta = se(3)$ can be represented by multiples of the generators

$$\delta = (\mathbf{u}, \omega)^T \in \mathbb{R}^6 \quad (14)$$

\mathbf{u} is the translation, $\mathbf{u} \in \mathbb{R}^3$

ω is the rotation (exactly the axis-angle representation, its normal is the rotation angle, and its direction is the rotation axis), $\omega \in \mathbb{R}^3$

There’s a 1-1 map between $SE(3)$ and $se(3)$

$$\begin{aligned} \delta &= \ln(C) \\ C &= \exp(\delta) \end{aligned} \quad (15)$$

B.2 Adjoint

consider left-multiplication and right-multiplication in $SE(3)$ group

$$C \cdot \exp(\delta) = \exp(\text{Adj}_C \cdot \delta) C \quad (16)$$

$$\text{Adj}_C = \begin{pmatrix} \mathbf{R} & t \times \mathbf{R} \\ \mathbf{0} & \mathbf{R} \end{pmatrix} \quad (17)$$

c.f. Skew-symmetric matrix of a vector $\alpha := [a, b, c]^T, \beta := [l, m, n]^T$

$$\alpha \times \beta = \begin{bmatrix} bn - cm \\ cl - an \\ am - bl \end{bmatrix} = \begin{bmatrix} 0 & -c & b \\ c & 0 & -a \\ -b & a & 0 \end{bmatrix} \begin{bmatrix} l \\ m \\ n \end{bmatrix} \quad (18)$$

So we define $\alpha_{\times}([\alpha]_{\times})$ as the skew-symmetric matrix in 18

B.3 Interpolation

two elements $a, b \in G$, we would like to interpolate between the two elements according to a parameter $t \in [0, 1]$, define an interpolation function

$$f : G \times G \times \mathbb{R} \rightarrow G \quad (19)$$

First define a group element that tasks a to b

$$d := b \cdot a^{-1}d \cdot a = b \quad (20)$$

Compute the corresponding Lie algebra vector and scale it by t

$$\mathbf{d}(t) = t \cdot \ln(d); d_t = \exp(\mathbf{d}(t)) \quad (21)$$

$$f(a, b, t) = d_t \cdot a$$

B.4 Gaussian Distribution

Consider a Lie group G and its Lie algebra vector space \mathfrak{g} , with k DoF. A mean transformation $\mu \in G$ and a covariance matrix $\Sigma \in \mathbb{R}^{k \times k}$. We can sample an element from the Gaussian distribution on G

$$\delta \sim \mathcal{N}(0; \Sigma)(\delta \in \mathfrak{g}) \quad (22)$$

$$x = \exp(\delta) \cdot \mu$$

c.f. Gaussian in \mathbb{R}^D

$$N(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) \right\} \quad (23)$$

C Proof of proposition 1

$$\begin{aligned} p(x^0|c) &= \int p(x^K|c)p(x^{0:K-1}|x^K, c)dx^{1:K} \\ &= \int p(x^K|c) \prod_{k=1}^K p(x^{k-1}|x^k, c)dx^{1:K} \\ &= \int p(x^K|c) \left(\prod_{k=n}^K p(x^{k-1}|x^k, c) \right) p(x^{n-2}|x^{n-1}, c) \left(\prod_{i=1}^{n-2} p(x^{i-1}|x^i, c) \right) dx^{1:K} \\ &= \int p(x^K|Tc) \left[\prod_{k=n}^K p_1(x^{k-1}|x^k, Tc) \right] p_2(Tx^{n-2}|x^{n-1}, Tc) \\ &\quad \left[\prod_{i=1}^{n-2} p_3(Tx^{i-1}|Tx^i, Tc) \right] dx^{1:K} \\ &= p(x^{n-1}|Tc) p(Tx^{n-2}|x^{n-1}, Tc) p(Tx^0|Tx^{n-2}, Tc) \\ &= p(Tx^0|Tc) \end{aligned} \quad (24)$$

D Proof proposition 2

As shown in 2 and 3, our backbones are theoretically $SE(3)$ equivariant and invariant, i.e.

$$\begin{aligned} \mathcal{E}_{inv}(O, \mathbf{H}_i^k; k, i) &= \mathcal{E}_{inv}(TO, \mathbf{H}_i^k; k, i) \\ T\mathcal{E}_{equiv}(O, \mathbf{H}_i^k; k, i) &= \mathcal{E}_{equiv}(TO, \mathbf{H}_i^k; k, i) \end{aligned} \quad (25)$$

For each element \mathbf{H}_i^K in A^K , we firstly use \mathcal{E}_{inv} to denoise $K-1$ steps, so for any $1 < k \leq K$, the denoise iteration satisfies

$$\hat{\mathbf{H}}_i^{k \rightarrow 0} = \mathcal{E}_{inv}(O, \mathbf{H}_i^k; k, i) \quad (26)$$

When the input observation is transformed by $SE(3)$ element, for the property of \mathcal{E}_{inv} , we have

$$\forall T \in SE(3) : \hat{\mathbf{H}}_i^{k \rightarrow 0} = \mathcal{E}_{inv}(TO, \mathbf{H}_i^k; k, i) \quad (27)$$

It means the predicted noise $\hat{\mathbf{H}}_i^{k \rightarrow 0}$ keeps invariant no matter what $SE(3)$ transformation is applied on the input observation. And then we carry the predicted noise into 8, it's obvious that \mathbf{H}_i^{k-1} is

also $SE(3)$ invariant. So we can infer that \mathbf{H}_i^1 is $SE(3)$ invariant to input observation. In terms of Markov transition, the first $K - 1$ transitions are p_1 -like.

$$p(\mathbf{H}_i^{k-1}|\mathbf{H}_i^k, O) = p(\mathbf{H}_i^{k-1}|\mathbf{H}_i^k, TO), 1 < k \leq K \quad (28)$$

For the last denoising iteration, we use a $SE(3)$ equivariant model to predict noise, so when the input observation is transformed, we have

$$\forall T \in SE(3) : T\hat{\mathbf{H}}_i^1 \rightarrow 0 = \mathcal{E}_{equiv}(TO, \mathbf{H}_i^1; 1, i) \quad (29)$$

Carry the result into 9, we will discover the final denoised action $\hat{\mathbf{H}}_i^0$ is $SE(3)$ Equivariant. In another word, the last Markov transition is p_2 -like.

$$p(\mathbf{H}_i^0|\mathbf{H}_i^1, O) = p(T\mathbf{H}_i^0|\mathbf{H}_i^1, TO) \quad (30)$$

Additionally, as the initial noisy action \mathbf{H}_i^K is sampled from Gaussian distribution, it's not conditioned on the observation.

$$p(\mathbf{H}_i^K|O) = p(\mathbf{H}_i^K|TO) \quad (31)$$

Combine 28, 30, 31 together and put them back into 24, we find the whole diffusion process for single action is $SE(3)$ equivariant. Joint all the $\mathbf{H}_i^0 (0 \leq i \leq T_p)$ into a sequence A^0 , it's easy to verify proposition 2 holds. In other word, **our predicted action sequence is theoretically $SE(3)$ equivariant to input observations.**

E Experiments showing p_1 -like transition is easier to learn

The properties of three different Markov transitions can be described as

$$\begin{aligned} p(y|x, c) &= p_1(y|x, Tc) \\ p(y|x, c) &= p_2(Ty|x, Tc) \\ p(y|x, c) &= p_3(Ty|Tx, Tc) \end{aligned} \quad (32)$$

In practice, we use the $SE(3)$ Transformer [17] with different input and output feature types to approximate the three types of transitions(Denoted as P1Net, P2Net and P3Net). In this validation experiment, we take a point cloud P with random orientation as observation (focusing solely on rotation for simplicity). The detailed input and output feature types are shown in table 1 . According to the features of $SE(3)$ Transformer, it's easy to verify the networks satisfy the corresponding equivariant properties.

Table 1: Input and Output Feature Types

	Input Feature	Output Feature	Supervision	Final Loss
P1Net	3 type-0	9 type-0	identity matrix	0.0002
P2Net	3 type-0	3 type-1	Pose of input pts	0.25
P3Net	1 type-1	3 type-1	Pose of input pts	0.27

For all three networks, the input feature consists of 3 scalar values attached to each point, and the output feature consists of 9 scalar values (after pooling across all points). For P1Net, the output is set as nine type-0 features, meaning the output remains invariant to the rotation of the input point cloud. We supervise the output by computing the L2 loss between it and a fixed rotation matrix. In contrast, for P2Net and P3Net, the output is treated as three type-1 features, which are supervised using the pose of the input point cloud. The only difference among the three networks is the input and output feature types, while all other hyperparameters remain the same.

After training for the same number of epochs, the loss curve of the three networks is shown in fig. 4. The experiments demonstrate that the invariant model (P1Net) is significantly easier to train compared to the equivariant models (P2Net and P3Net), as it is expected to output the same values regardless of the transformation applied to the input point cloud. Additionally, we observe that the

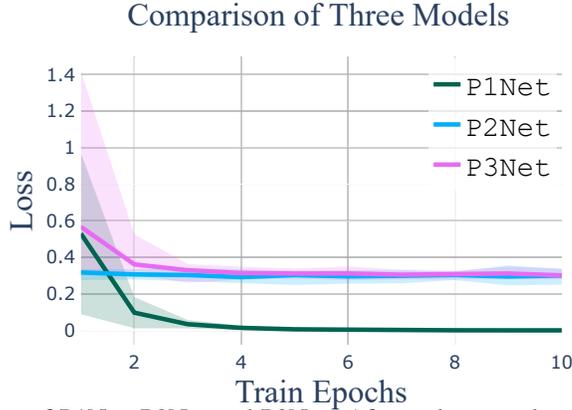


Figure 4: Loss curve of P1Net, P2Net and P3Net. After only several gradient descent, the loss of P1Net converges almost to 0, while the losses of P2Net and P3Net do not decrease obviously.

use of higher-type features in P2Net and P3Net results in increased memory requirements and longer inference times.

Therefore, in ET-SEED, the diffusion process consists of $K - 1$ p_1 -like denoising steps and only one p_2 -like step.

F Implementation of $SE(3)$ equivariant and invariant backbones

Here we introduce the implementation of our true $SE(3)$ equivariant backbone \mathcal{E}_{equiv} and invariant backbone \mathcal{E}_{inv} $SE(3)$ Transformer [17].

In general, each module consists of 2 $SE(3)$ Transformers, called as *pos_net* and *ori_net*, outputting translation and rotation separately. As the output of $SE(3)$ Transformer is per-point features, we implement a mean pooling over all points to get global features.

F.1 Invariant module

We set the output of *ori_net* as 6 type0 features, and then implement Schimidt orthogonalization to get rotation matrix. As the 6 type0 features are $SE(3)$ invariant, the rotation matrix is also invariant.

We set the output of *pos_net* as 3 type0 features, which is naturally invariant to any $SE(3)$ transformation of the point cloud, guaranteed by the translation invariance of $SE(3)$ Transformer.

Finally we combine the translation and rotation parts to a 4×4 matrix, and it is invariant to any $SE(3)$ transformations of input point cloud.

F.2 Equivariant module

We set the output of *ori_net* as 2 type1 features, and then implement Schimidt orthogonalization to get rotation matrix. As the 2 type1 features are $SE(3)$ equivariant, the rotation matrix is also equivariant.

We set the output of *pos_net* as 2 type1 feature and 3 type0 feature(denoted as *offset*). First we implement Schimidt orthogonalization on the 2 type1 feature, get a rotation matrix (denoted as R). Additionally, we denote the mass center of the input point cloud as $\mathcal{M} := \frac{1}{N} \sum_{i=1}^N x_i$, x_i is the coordinate of the i -th point. Then we can write the predicted translation t as

$$t(\mathcal{X}) = \mathcal{M} + R \cdot offset \tag{33}$$

We can prove this translation vector is equivariant to any $SE(3)$ transformation of the input point-cloud \mathcal{X} . When the input point cloud is transformed, $\mathcal{X}' = R_{data}\mathcal{X} + t_{data}$.

$$\begin{aligned} t'(\mathcal{X}) &= (R_{data}\mathcal{M} + t_{data}) + R_{data}R \cdot offset \\ &= R_{data}(\mathcal{M} + R \cdot offset) + t_{data} \\ &= R_{data}t(\mathcal{X}) + t_{data} \end{aligned} \quad (34)$$

Finally we combine the translation and rotation parts to a 4×4 matrix, and it is equivariant to any $SE(3)$ transformations of input point cloud.

G Equivariant Diffusion on $SE(3)$ manifold

The noisy action of step k can be represented as

$$\mathbf{H}^k = \underbrace{\text{Exp}(\gamma\sqrt{1-\bar{\alpha}_t}\varepsilon)}_{\text{Perturbation}} \underbrace{\mathcal{F}(\sqrt{\bar{\alpha}_t}; \mathbf{H}^0, \mathbb{H})}_{\text{Interpolation}}, \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (35)$$

The first term, $\text{Exp}(\gamma\sqrt{1-\bar{\alpha}_t}\varepsilon)$ is a random noise on $SE(3)$ manifold, aiming to randomize the diffusion process. According to the Gaussian distribution on $SE(3)$ (appendix B.4), a $SE(3)$ Gaussian variable can be written as the Exp of a Gaussian variable on $\mathfrak{se}(3)$. So we first randomly sample a 6D noise $\varepsilon \in \mathbb{R}^6$ from unit Gaussian distribution, then scale it by a scheduler factor $\gamma\sqrt{1-\hat{\alpha}_t}$ to control the magnitude of the perturbation at different steps. Finally we use the Exp map to convert the variable back to $SE(3)$.

The second term is an interpolation on $SE(3)$ manifold between \mathbf{H}^0 and \mathbb{H} . The idea behind this function is, first project the $SE(3)$ transformation to $\mathfrak{se}(3)$, perform linear interpolation in this tangent space, and then convert the interpolated vector back to $SE(3)$ to obtain the interpolated transformation. One can refer to Jiang et al. [16] for more details. Formally, the interpolation function \mathcal{F} can be expressed as

$$\mathcal{F}(\sqrt{\bar{\alpha}_t}; \mathbf{H}_0, \mathbb{H}) = \text{Exp}((1-\sqrt{\bar{\alpha}_t}) \cdot \log(\mathbb{H}\mathbf{H}_0^{-1})) \mathbf{H}_0 \quad (36)$$

H Simulation Experiments– Further Details

Table 2: Success rates (\uparrow) and standard deviation of different tasks in simulation.

Method	Open Bottle Cap				Open Door				Rotate Triangle			
	T		NP		T		NP		T		NP	
	25	50	25	50	25	50	25	50	25	50	25	50
DP3 [32]	65±4.5	76±5.5	11±4.2	14±6.5	61±2.24	72±2.74	9±3.54	16±5.48	67±2.74	89±2.24	5±2.24	10±2.74
DP3+Aug	35.0±5.0	44±4.2	38±4.47	46±7.42	43±8.37	54±6.52	30±4.18	40±8.22	35±3.54	42±4.47	32±5.70	41±4.18
EquiBot [14]	63±2.74	73±2.74	63±5.70	77±7.58	56±2.24	72±2.24	58±7.58	77±7.58	67±2.74	84±2.24	64±8.66	86±5.48
ET-SEED(Ours)	67±2.74	81±2.24	74±6.52	82±2.74	66±2.24	75±2.74	66±2.74	76±2.24	83±2.24	93±2.74	85±2.24	89±4.18
Method	Calligraphy				Fold Garment				Fling Garment			
	T		NP		T		NP		T		NP	
	25	50	25	50	25	50	25	50	25	50	25	50
DP3 [32]	28±2.74	50±3.54	0±0.00	3±2.74	44±2.24	60±4.18	4±5.48	8±4.48	36±5.48	67±4.48	4±5.48	9±8.22
DP3+Aug	8±2.74	21±4.18	3±2.24	12±11.51	13±5.70	27±7.58	17±10.37	31±9.62	28±7.58	38±4.48	11±4.18	31±2.24
EquiBot [14]	24±5.48	43±8.37	14±10.84	40±10.61	34±4.18	58±2.74	33±2.74	60±7.90	35±6.12	61±6.52	36±6.52	64±8.22
ET-SEED(Ours)	38±2.74	55±3.54	36±6.52	54±8.22	47±2.74	67±2.74	49±2.24	69±4.18	50±5.00	67±4.48	48±4.47	62±5.70

Table 3: $SE(3)$ Geodesic distances (\downarrow) of different tasks in simulation.

Method	Open bottle cap				Open Door				Rotate Triangle			
	T		NP		T		NP		T		NP	
	25	50	25	50	25	50	25	50	25	50	25	50
DP3 [32]	0.257	0.197	1.413	1.785	0.384	0.354	0.478	0.442	0.265	0.192	1.812	1.627
DP3+Aug	0.283	0.234	0.276	0.218	0.391	0.315	0.442	0.329	0.247	0.187	0.578	0.447
EquiBot [14]	0.194	0.151	0.197	0.170	0.241	0.224	0.266	0.228	0.197	0.107	0.214	0.099
ET-SEED (Ours)	0.133	0.114	0.127	0.124	0.127	0.101	0.121	0.128	0.098	0.082	0.104	0.087

Method	Calligraphy				Fold Garment				Fling Garment			
	T		NP		T		NP		T		NP	
	25	50	25	50	25	50	25	50	25	50	25	50
DP3 [32]	0.305	0.241	4.988	4.662	0.479	0.298	4.466	4.179	0.529	0.348	4.993	4.365
DP3+Aug	0.354	0.337	4.752	4.365	1.318	0.976	1.524	1.219	1.318	0.976	1.524	1.219
EquiBot [14]	0.291	0.117	0.282	0.129	0.368	0.293	0.387	0.288	0.418	0.343	0.437	0.338
ET-SEED (Ours)	0.124	0.083	0.121	0.089	0.299	0.149	0.287	0.136	0.349	0.179	0.337	0.186

H.1 Task Setting:

- *Rotate Triangle*: A robotic arm with 2D anchor pushes the triangle to a target 6D pose. The task reward is computed as the percentage of the Triangle shape that overlaps with the target Triangle pose.
- *Open Bottle Cap*: A bottle with a cap is placed at a random position in Workspace, and a robot arm is tasked with opening the cap. In this task, the demonstrations show robots Unscrewing bottle cap with parallel gripper. Success in this task depends on whether the bottle cap is successfully opened. Note that, due to simulator constraints, opening the bottle cap simply involves lifting it upward without the need to twist it first.
- *Open Door*: This task evaluates the manipulation of articulated objects. The model is required to generate trajectories to open doors positioned at various orientations. The demonstration is given as: We initialize the gripper at a point p sampled on the handle of the door and set the forward orientation along the negative direction of the surface normal at p . And then we pull the door by a degree. Different from door pushing, we perform a grasping at contact point p for pulling. Success in this task is determined by the opening angle of the door.
- *Robot Calligraphy*: This long-horizon task involves using a robot arm to write complex Chinese characters on paper, accounting for different orientations. Success in this task is determined by the aesthetic quality and accuracy of the characters or patterns formed, which should closely resemble the target trajectory.
- *Fold Garment*: A long-horizon task involving deformable object manipulation, where a robot folds a long-sleeved garment. The robot folds the sleeves inward along the garment’s central axis, then gathers the lower edge of the garment and folds it upward, aligning it with the underarm region. A folding succeeds when the Intersection-over-Union (IOU) between the target and the folded garments exceeds a bar [34, 35, 36, 37].
- *Fling Garment*: A dual-arm task for manipulating deformable objects. The robot grasps the two shoulder sections of a wrinkled dress, lifting it to allow the fabric is clear of the surface. Then flings the garment it to flatten the fabric, and then places it back onto a flat surface. Success is determined by the projection area of the flattened garment.

Some of the six tasks are exactly $SE(3)$ equivariant, and some are partially. In the *Open Door* task, the manipulation trajectory is exactly equivariant with the point cloud of the door. In the *Robot Calligraphy* task, the manipulation trajectory is exactly equivariant with the point cloud of the paper and the handwriting that has been written. In the *Open Bottle Cap* task, the manipulation trajectory is exactly equivariant with the point cloud of the bottle. In the *Rotate Triangle* task, as we always add same transformation on initial pose and target pose of the triangle, the manipulation trajectory

is exactly equivariant with the point cloud of the triangle. In the *Fling Garment* and *Fold Garment* tasks, the trajectories are not exactly equivariant with the initial observation of the garment, as the deformation differs in each initialization. But even so, our method can also outperform baselines.

The intuition is, by ensuring end-to-end $SE(3)$ equivariance, our model treats a point cloud in the observation space and all the point clouds possible by $SE(3)$ transformation as an equivalence class, so the observation space is reduced to the quotient space of the observation space over the $SE(3)$ group. Once the $SE(3)$ equivariance is naturally ensured, the network can focus on the geometric features of the object, so it still has the ability to generalize to the deformation and geometry of objects.

H.2 Qualitative Results.

Table 2 and 3 provide a quantitative comparison between our method and the baseline. Both DP3 and its augmented variant demonstrate strong performance in the training setting (T), but they exhibit a significant drop in performance when faced with New Poses (NP) scenarios. This highlights that merely incorporating data augmentation is insufficient for the model to generalize effectively to unseen poses. Instead, leveraging equivariance proves essential for enhancing spatial generalization.

While EquiBot achieves commendable results in both success rate and \mathcal{D}_{geo} , it struggles with more complex, long-horizon tasks such as Calligraphy and Fold Garment. Also, when less demonstrations are given, the performance is not satisfactory. These challenges stem from the inherent complexity of its diffusion process design, where maintaining equivariance in each Markov transition adds substantial difficulty to the learning task.

In contrast, ET-SEED consistently outperforms across all six tasks, with minimal performance drop when facing unseen object poses. This advantage is especially pronounced when using a limited number of demonstrations, showcasing T-SEED’s superior data efficiency, manipulation proficiency, and spatial generalization ability.

I Real World Experiments– Further Details



Figure 5: Visualizations of the real-world environments used in our experiments. The tasks are performed using multiple Microsoft Azure Kinect cameras and Intel® RealSense for point cloud fusion and a Panda robotic arm for execution.

Table 4: Success rates in real-world robot experiments.

Method	Open Bottle Cap	Open Door	Calligraphy	Fold Garment
DP3	0.2	0.2	0.0	0.1
DP3+Aug	0.2	0.3	0.0	0.2
EquiBot	0.6	0.5	0.0	0.3
ET-SEED (Ours)	0.8	0.6	0.4	0.6

Task Description. In the task of **Open Bottle Cap**, unlike in a simulator, the process in the real world involves first twisting the cap and then lifting it off to open. The bottle is initially placed at a random position on the table. For **Opening Door**, the initial position of the cabinet is randomly

determined. In **Calligraphy**, we use flat-bristled brushes and watercolor paper, which is randomly positioned on the table. In **Cloth Folding**, the limited workspace of the Franka robot means it cannot reach every possible location. Therefore, the placement of clothes is not random but is instead based on locations accessible to the robot.

Qualitative Results. Results for our real robot tasks are given in Table 4. Consistent with our simulation findings, in real world experiments, ET-SEED performs better than baselines in all the four tasks, given only 20 demonstrations. The evaluation shows the effectiveness and spatial generalization ability of our method.

J Ablation Study

Table 5: Ablation studies.

Design	Average
Ours w/o SE(3)	24±4.48
Ours w/o Eqv-Diff	57±6.52
Ours	76±2.24

Ablation Studies. We conduct ablation studies on the New Pose (NP) scenario of the representative *Opening Door* task to evaluate the effectiveness of different components of our approach:

- **Ours w/o SE(3):** Our method without $SE(3)$ invariance and equivariance in the backbone architecture. In this variant, we use a standard PointNet++ to predict noise at each step.
- **Ours w/o Eqv-Diff:** Our method without the $SE(3)$ equivariant denoising process. Instead, we use a non-equivariant diffusion process (DDIM), following the approach of Ze et al. [32].

Table 5 shows quantitative comparisons with ablations. Clearly each component improves our method’s capability.