

E2PN: Efficient SE(3)-Equivariant Point Network

Minghan Zhu
University of Michigan
minghanz@umich.edu

Maani Ghaffari
University of Michigan
maanigj@umich.edu

William A Clark
Cornell University
wac76@cornell.edu

Huei Peng
University of Michigan
hpeng@umich.edu

Abstract

This paper proposes a convolution structure for learning SE(3)-equivariant features from 3D point clouds. It can be viewed as an equivariant version of kernel point convolutions (KPCov), a widely used convolution form to process point cloud data. Compared with existing equivariant networks, our design is simple, lightweight, fast, and easy to be integrated with existing task-specific point cloud learning pipelines. We achieve these desirable properties by combining group convolutions and quotient representations. Specifically, we discretize $SO(3)$ to finite groups for their simplicity while using $SO(2)$ as the stabilizer subgroup to form spherical quotient feature fields to save computations. We also propose a permutation layer to recover $SO(3)$ features from spherical features to preserve the capacity to distinguish rotations. Experiments show that our method achieves comparable or superior performance in various tasks, including object classification, pose estimation, and keypoint-matching, while consuming much less memory and running faster than existing work. The proposed method can foster the development of equivariant models for real-world applications based on point clouds.

1. Introduction

Processing 3D data has become a vital task today as demands for automated robots and augmented reality technologies emerge. In the past decade, computer vision has significantly succeeded in image processing, but learning from 3D data such as point clouds is still challenging. An important reason is that 3D data presents more variations than 2D images in several aspects. For example, the rigid body transformations in 2D only have 3 degrees of freedom (DoF) with 1 for rotations. In 3D space, the DoF is 6, with 3 for rotations. The 2D translation equivariance is a key factor in the success of convolutional neural networks (CNNs)

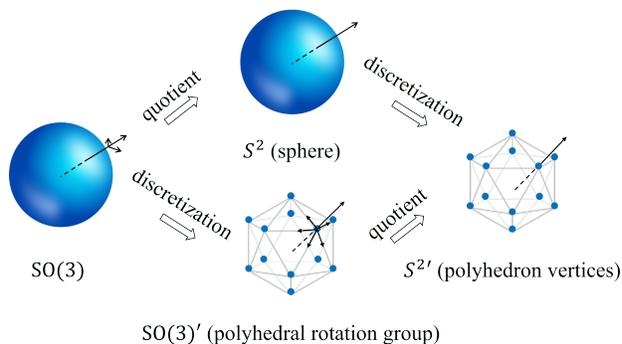


Figure 1. Our method achieves higher efficiency by working with smaller feature maps defined on $S^{2'} \times \mathbb{R}^3$ rather than $SO(3)' \times \mathbb{R}^3$ (' denotes discretization). \mathbb{R}^3 is omitted in the figure. The black arrows in each space represent elements. The top and bottom paths are equivalent, showing the relations among different representations.

in image processing, but it is not enough for 3D tasks.

Generally speaking, *equivariance* is a property for a map such that given a transformation in the input, the output changes in a predictable way determined by the input transformation. It drastically improves generalization as the variance caused by the transformations is captured via the network by design. Take CNNs as an example, the equivariance property refers to the fact that a translation in the input image results in the same translation in the feature map output from a convolution layer. However, conventional convolutions are not equivariant to rotations, which becomes problematic, especially when we deal with 3D data where many rotational variations occur. In response, on the one hand, data augmentations with 3D rotations are frequently used. On the other hand, equivariant feature learning emerges as a research area, aiming to generalize the translational equivariance to broader transformations.

A lot of progress has been made in group-equivariant feature learning. The term *group* encompasses the 3D rotations and translations, which is called the special Euclidean

group of dimension 3, denoted $SE(3)$, and also other more general types of transformations that represent certain symmetries. While many methods have been proposed, equivariant feature learning has not yet become the default strategy for 3D deep learning tasks. From our understanding, two major reasons hinder the broader application of equivariant methods. First, networks dealing with continuous groups typically require specially designed operations not commonly used in neural networks, such as generalized Fourier transform and Monte Carlo sampling. Thus, incorporating them into general neural networks for 3D learning tasks is challenging. Second, for the strategy of working on discretized (finite) groups [6, 19], while the network structures are simpler and closer to conventional networks, they usually suffer from the high dimensionality of the feature maps and convolutions, which causes much larger memory usage and computational load, limiting their practical use.

This work proposes E2PN, a convolution structure for processing 3D point clouds. Our proposed approach can enable $SE(3)$ -equivariance on any network with the KPConv [35] backbone by swapping KPConv with E2PN. The equivariance is up to a discretization on $SO(3)$. We leverage a quotient representation to save computational and memory costs by reducing $SE(3)$ feature maps to feature maps defined on $S^2 \times \mathbb{R}^3$ (where S^2 stands for the 2-sphere). Nevertheless, we can recover the full 6 DoF information through a final permutation layer. As a result, our proposed network is $SE(3)$ -equivariant and computationally efficient, ready for practical point-cloud learning applications.

Overall, this work has the following contributions:

- We propose an efficient $SE(3)$ -equivariant convolution structure for 3D point clouds.
- We design a permutation layer to recover the full $SE(3)$ information from its quotient space.
- We achieve comparable or better performance with significantly reduced computational cost than existing equivariant models.
- Our implementation is open-sourced at <https://github.com/minghanz/E2PN>.

Readers can find preliminary introductions to some related background concepts in the appendix.

2. Related work

Group convolutions: In 2016, Cohen and Welling proposed G-CNN [6], enabling equivariance beyond translations on 2D images with generalized G-convolutions (group convolutions) over the group of 90-degree rotations, which is one of the earliest efforts in equivariant deep learning. Group convolution is similar to conventional convolutions but has an extended domain for feature maps and kernels.

The idea was then applied to different networks to enable equivariance for $SE(2)$ [6, 19], $SO(3)$ [5], $SE(3)$ [2, 42], and $E(3)$ [41] groups up to some discretization. The idea mainly works with finite (discretized) groups, as it is convenient to parameterize feature maps and kernels on the discretized group elements just as on pixel grids. Group convolutions have a relatively simple structure, making them straightforward to apply, but a major downside is that the lifted domain of features and kernels causes higher computational and memory costs, and the problem is more prominent when the group is large. Groups convs can also work with continuous groups, for example, with the help of Monte Carlo (MC) estimation as in [13], but they suffer from a large memory burden in MC sampling when the number of layers grows [30].

Steerable CNNs: Another line of work is steerable CNNs. Instead of augmenting the domain of feature maps, steerable CNNs generalize the space of feature values to be *steerable*, i.e., the feature values transform predictably as the input transforms. The way the feature transforms is called the group representation in the feature space, governed by the feature *type*. Features of *scalar* type are kept unchanged under group actions, and we call the group representation trivial. For *vector* or *tensor* feature types, the representation is not trivial, and the features will change with group actions, for example, through rotation matrix multiplication. Steerable CNNs work for both discretized and continuous groups. One may freely design the feature types based on pre-determined basic types and corresponding representations (irreducible representations, i.e., *irreps*) as building blocks for a given group. Group convolutions can be viewed as a special case of steerable CNNs with *regular* representations, i.e., channels of a feature vector undergo permutations when transformed. Examples of steerable CNNs include [4, 9, 36, 40, 43]. While the framework of steerable CNNs generalizes group convolutions with more flexibility, it requires a good understanding of representation theory and involves generalized Fourier analysis when working with a continuous group, which makes the structure complicated and challenging to apply in real-world applications broadly. The work of [21, 39] also found empirically that steerable CNNs with irreps underperform regular representations in certain tasks.

Theoretical progress: There has been a lot of progress in the theoretical development of equivariant networks that are not group-specific. [24] developed a general formulation for steerable CNNs with scalar type features, and [7, 8] generalized it to arbitrary feature types. It is proven that all equivariant linear maps can be written as convolutions. The formulation from the perspective of Fourier analysis was presented in [46]. [39] found the general form of solution for the equivariant kernels for $E(2)$ group, later generalized to any compact group [26]. [14] proposed an algorithm to solve for the equivariance constraint for arbitrary

matrix group. The formulation of group convolution based on MC estimation was proposed for any Lie group with [13] or without [30] surjective exponential maps. Equivariant non-linear layers like transformers [15, 21] and equivariant set and graph networks [14, 22, 32] are also proposed, but they are not the focus of this paper.

Applications of equivariant learning in perception tasks: We want to highlight a few equivariant networks that gain attention in perception applications due to their simplicity and practicality. Vector Neurons [10] is a PointNet-like $SO(3)$ -equivariant network for 3D point cloud learning, later applied to point cloud registration [49] and manipulation [34]. It can be viewed as a special case of TFN [36] with type-1 features and self-interactions only. E2CNN [39] as an $SE(2)$ -equivariant framework was applied in several image processing tasks [27, 33], given its generality and user-friendly library. DEVIANT [25] applied scale-equivariant convolutions in monocular 3D object detection. EPN [2] is a group convolution network with $SE(3)$ -equivariance for 3D point cloud learning based on KPConv [35] and was used in practice, for example in place recognition task [29]. We also position this proposed work in this category, aiming to promote the application of equivariant learning with our efficient and easy-to-use design. Our work is developed based on EPN, which also serves as a major baseline in this paper.

3. Methodology

3.1. Overview of the idea

We first explain why our proposed method is more efficient. Then we discuss how we gain efficiency without sacrificing expressiveness.

3.1.1 Improved efficiency with quotient features

Following the idea of group convolutions, one needs to extend the domain of feature maps from the Euclidean space to the group space, from \mathbb{R}^3 to $SE(3) \cong SO(3) \times \mathbb{R}^3$ in our case, where we want $SE(3)$ -equivariant features for 3D point clouds. We then discretize $SO(3)$ to a finite group denoted as $SO(3)'$ (we use $'$ to denote discretization in this paper). The discretization of \mathbb{R}^3 is taken care of by KPConv, which is the non-equivariant counterpart of our method, thus, not discussed here. Up to now, we can obtain a $SE(3)$ -equivariant version of KPConv [35] realized through group convolution, which is what EPN [2] presents.

The specific form of the finite group $SO(3)'$ is introduced later in Sec. 3.2, but it can be geometrically understood as the set of all rotations that keep a Platonic solid (i.e., a convex and regular 3D polyhedron) unchanged. As illustrated at the bottom of Fig. 1, the set of rotations can be enumerated by counting the number of vertices, repre-

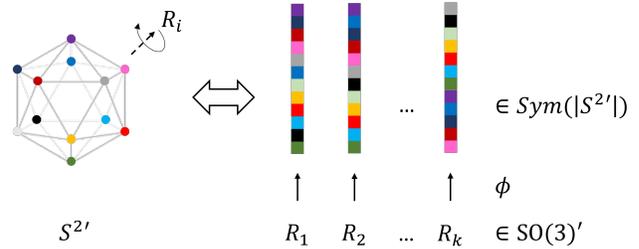


Figure 2. Illustration of recovering $SO(3)'$ features from $S^{2'}$ features through permutations. The left and right subplots are from the geometric and algebraic views. See Sec. 3.1.2 for more explanation.

sending rotations taking a given vertex to any vertices, multiplied by the number of edges connected to a single vertex, representing the rotations that keep a vertex fixed. The same result can be achieved by counting the faces or the edges, but we stick with vertices in the following discussion.

In comparison, we propose to define feature maps on $S^2 \times \mathbb{R}^3$, where S^2 is the sphere space. We have $S^2 = SO(3)/SO(2)$, meaning that it is the *quotient space* of $SO(3)$ given the stabilizer subgroup $SO(2)$. To understand the quotient space intuitively, we can see that all rotations can be grouped by the destination of a point on the sphere (e.g., the north pole) after the rotation. All rotations bringing the north pole to the same destination point form a *coset*. They are related by rotations around the axis passing through that point, forming a subgroup isomorphic to $SO(2)$. Thus S^2 is the quotient of $SO(3)$ and $SO(2)$. In the discretized setup, as depicted in Fig. 1, $SO(2)'$ is discretized by the number of edges connected to a single vertex. The quotient $S^{2'} = SO(3)'/SO(2)'$ corresponds to the vertices on the Platonic solid.

In general, $|S^{2'}| = |SO(3)'|/|SO(2)'| < |SO(3)'|$, thus the size of feature maps and kernels defined on $S^{2'} \times \mathbb{R}^3$ is much smaller than those on $SO(3)' \times \mathbb{R}^3$. The convolution operation on the former space also requires smaller costs. It is the major reason our method is much more efficient than EPN. There is another design, called *symmetric kernel*, which further improves the efficiency of our method, but we will refer to Sec. 3.2.3 for details.

3.1.2 Information recovery with permutations

With the quotient feature map, all rotations moving the north pole to the same point on a sphere are represented by the same point. Thus we immediately lose the ability to distinguish among these rotations, which is a problem if our task is to learn the pose, for example, in the point cloud registration tasks.

However, with our proposed permutation layer, we can distinguish every element in $SO(3)'$ from the feature maps on $S^{2'}$. The key observation is that the action of $SO(3)$ on S^2 is *faithful*, i.e., $\forall R \in SO(3)$ and $R \neq I$, $\exists x \in S^2$,

s.t. $Rx \neq x$. It means that an action of $SO(3)$ other than identity will always cause a change on the S^2 feature map when looking at all points on S^2 simultaneously. In the discretized setup, it means that there is an injective map $\phi : SO(3)' \rightarrow Sym(|S^{2'}|)$, where $Sym(|S^{2'}|)$ stands for the symmetric group, i.e., the collection of all permutations of a set of size $|S^{2'}|$. An element in $Sym(n)$ is a bijective map $\{1, \dots, n\} \rightarrow \{1, \dots, n\}$, permuting the indices. In other words, each rotation corresponds to a unique permutation of the $S^{2'}$ elements, from which we can distinguish each rotation from the feature map defined on $S^{2'}$.

Here we provide the specific form of the permutation layer. Given ϕ mentioned above and a feature map $f : S^{2'} \rightarrow \mathbb{R}^m$, we can build a feature map $\tilde{f} : SO(3)' \rightarrow \mathbb{R}^{mn}$ defined as:

$$\tilde{f}(R) = [f(x[\phi_R(1)]), f(x[\phi_R(2)]), \dots, f(x[\phi_R(n)])] \quad (1)$$

where $R \in SO(3)'$, $n = |S^{2'}|$, and $x[i] \in S^{2'}$ for $i = 1, \dots, n$. Given that ϕ is injective, $\tilde{f}(R_1) \neq \tilde{f}(R_2)$ when $R_1 \neq R_2$. Thus we can distinguish $SO(3)'$ rotations from \tilde{f} . See Fig. 2 for an illustration.

3.2. Specific form of convolution

3.2.1 Recap of KPConv

A conventional 3D convolution can be written as:

$$[\kappa * f](x) = \int_{\mathbb{R}^3} \kappa(t) f(x+t) dt = \sum_{t \in \mathbb{R}^{3'}} \kappa(t) f(x+t), \quad (2)$$

where $x \in \mathbb{R}^3$, and the right hand side is after discretization. We use *correlations* to implement convolutions in this paper following conventions in deep learning. For processing point clouds, Eq. (2) could be tricky in implementation because it is challenging to align κ with f when there is no grid. The strategy of KPConv is to have a set of kernel points for κ and to gather features to the kernel point coordinates from the input points where f is defined so that they are aligned before the convolution. It is depicted in Fig. 3 (a), and we need to replace f with \hat{f} in Eq. (2) where $\hat{f}(x) = \sum_{y \in \mathcal{N}_x} w(|y-x|) f(y)$ and w is a scalar weight function based on distance. \hat{f} is the features gathered from neighboring input points to align with the kernel points. In this case, the input and output feature map f and $[\kappa * f]$ are defined on \mathbb{R}^3 , while the kernel κ is defined on $\mathbb{R}^{3'}$, i.e., coordinates of the set of kernel points. The gathered feature \hat{f} when calculating the convolution at $x \in \mathbb{R}^3$ is defined at $x + \mathbb{R}^{3'}$. Notice that we do not consider the deformable mode of KPConv in this paper.

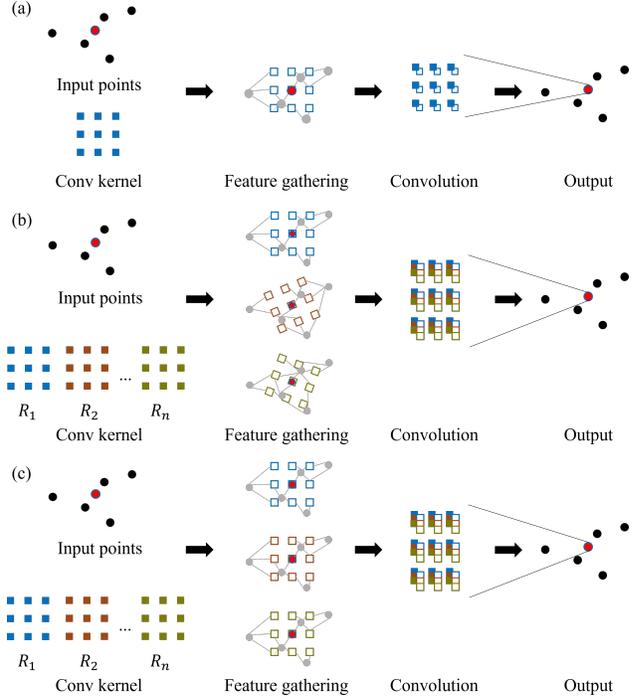


Figure 3. Illustration of different versions of KPConv. (a) original KPConv (see Sec. 3.2.1). (b) KPConv on finite group or quotient space without the symmetric kernel (see Sec. 3.2.2). (c) KPConv on group or quotient space with the symmetric kernel (see Sec. 3.2.3). (b,c) compared with (a): the convolution kernel is larger because it is defined on a higher dimension. (c) compared with (b): the feature gathering is more efficient because the kernel points are symmetric to the rotations in the rotation group.

3.2.2 KPConv on group and on quotient space

To extend KPConv to a group convolution [6] with finite group G' , we modify Eq. (2) to

$$[\kappa * f](g) = \sum_{g_t \in G'} \kappa(g_t) \hat{f}(g \cdot g_t), \quad (3)$$

where in our case, $G' = SE(3)' = SO(3)' \times \mathbb{R}^{3'}$ is where κ is defined, $g \in G^\dagger = SO(3)' \times \mathbb{R}^3$ is where the input and output feature map is defined. \hat{f} is defined on $g \cdot G'$ when calculating the convolution at g . Denote an element $g \in SE(3)$ as (R, t) where $R \in SO(3)$ and $t \in \mathbb{R}^3$, the binary operation \cdot for $SE(3)$ can be specified as $(R_1, t_1) \cdot (R_2, t_2) = (R_1 R_2, t_1 + R_1 t_2)$, same for the discretized case. EPN [2] follows this form of convolution. However, as discussed in Sec. 3.1.1, the size of group $|SE(3)'| = |SO(3)'| |\mathbb{R}^{3'}|$, making it expensive to store the feature map and compute the convolution. To alleviate this issue, EPN has to conduct convolutions in $SO(3)'$ and $\mathbb{R}^{3'}$ separately (i.e., separable convolutions [2, 3]), so that the computational cost is reduced.

We use quotient feature maps to address this issue. To conduct convolutions on the quotient space $X =$

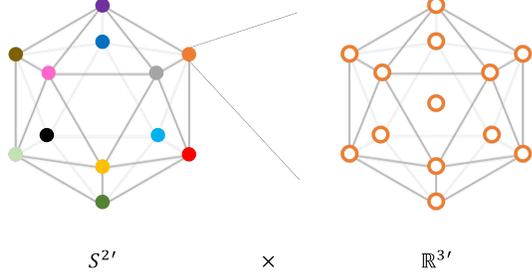


Figure 4. Visualization of the symmetric kernel κ used in our work, where $\mathbb{R}^{3'} = \{rS^{2'} \cup 0 \mid r > 0\}$.

$SE(3)/SO(2) = S^2 \times \mathbb{R}^3$, we cannot directly use Eq. (3), because the \cdot operation is not defined between two elements in X since X is not a group. It is pointed out in [7, 8] that we can write convolutions on the quotient space as

$$[\kappa * f](x) = \sum_{x_t \in X'} \kappa(x_t) \hat{f}(s(x) \cdot x_t), \quad (4)$$

where $X' = S^{2'} \times \mathbb{R}^{3'}$ is the domain of κ , $X^\dagger = S^{2'} \times \mathbb{R}^3$ is the domain of f and $[\kappa * f]$, and \hat{f} is defined on $s(x) \cdot X'$. $s(x)$ is called the *section map*, $s : X' \rightarrow G'$ (or $X \rightarrow G$ in the continuous case), mapping $x \in X'$ to an element of G' in the corresponding coset. With the section map, the \cdot operation denotes the action of the group G on quotient space X , i.e., $\cdot : G \times X \rightarrow X$. Denote an element in $S^2 \times \mathbb{R}^3$ as $(R\mathbf{n}, t)$, where $R \in SO(3)$, \mathbf{n} is the north pole point on the unit sphere $(0, 0, 1)$. Then $R\mathbf{n}$ represents arbitrary points on the sphere S^2 . The action \cdot can then be written as $(R_1, t_1) \cdot (R_2\mathbf{n}, t_2) = (R_1R_2\mathbf{n}, t_1 + R_1t_2)$.

3.2.3 Symmetric kernels

Now we introduce the specific form of $\mathbb{R}^{3'}$, i.e., the location of kernel points in KPConv. In our design, $\mathbb{R}^{3'} = \{rS^{2'} \cup 0 \mid r > 0\}$, i.e., the set of vertices of the Platonic solid with radius r and the origin point. $\mathbb{R}^{3'}$ is very similar to $S^{2'}$, because we want to make the kernel symmetric to $SO(3)'$. More precisely, we desire $\mathbb{R}^{3'}$ to be closed under the action of $SO(3)'$. There are two reasons as follows.

Steerability constraint To make the convolution on quotient space equivariant, defining a valid form of convolution as Eq. (4) is not the whole story. The kernel values must satisfy a condition called the *steerability constraint*, which is required for all steerable CNNs. More background knowledge can be found in [8]. In our case, the steerability constraint is

$$\kappa(x) = \kappa(R_z \cdot x), \forall x \in X', \forall R_z \in SO(2)', \quad (5)$$

where R_z is a z -axis rotation. The derivation of Eq. (5) from the general form of steerability constraints can be found in

Rotation groups	Tetrahedral T	Octahedral O		Icosahedral I	
Platonic solids	 Tetrahedron	 Cube	 Octahedron	 Dodecahedron	 Icosahedron
Size	$ SO(3)' = 12$ $ S^{2'} = 4$ $ SO(2)' = 3$	$ SO(3)' = 24$ $ S^{2'} = 8$ $ SO(2)' = 3$	$ SO(3)' = 24$ $ S^{2'} = 6$ $ SO(2)' = 4$	$ SO(3)' = 60$ $ S^{2'} = 20$ $ SO(2)' = 3$	$ SO(3)' = 60$ $ S^{2'} = 12$ $ SO(2)' = 5$

Figure 5. Illustration of all five types of Platonic solids and their corresponding finite rotation groups.

the appendix. Here the \cdot operation inherits from the action of G' on X' since $SO(2)' \subset SO(3)' \subset G'$. Specifically, we have $R_z \cdot (R\mathbf{n}, t) = (R_z, 0) \cdot (R\mathbf{n}, t) = (R_zR\mathbf{n}, R_zt)$. Replace x with $(R\mathbf{n}, t)$ in Eq. (5), and we have

$$\kappa(R\mathbf{n}, t) = \kappa(R_zR\mathbf{n}, R_zt), \quad (6)$$

$\forall R\mathbf{n} \in S^{2'}, \forall t \in \mathbb{R}^{3'}, \forall R_z \in SO(2)'$. Notice that it implies that we need $R_zt \in \mathbb{R}^{3'}, \forall t \in \mathbb{R}^{3'}, \forall R_z \in SO(2)'$, so that κ is defined on the right hand side of Eq. (6). In other words, the kernel points $\mathbb{R}^{3'}$ need to be closed under $SO(2)'$. Obviously, having $\mathbb{R}^{3'}$ closed under $SO(3)'$ is a sufficient condition for this.

For the $S^{2'}$ dimension, since $S^{2'}$ is symmetric (closed) to $SO(3)'$ by definition and $SO(2)' \subset SO(3)'$, we always have $R_zR\mathbf{n} \in S^{2'}$ and κ is always defined.

Efficient feature gathering Another important reason for the design choice of $SO(3)'$ -symmetric kernels is that it enables more efficient feature gathering. Consider a spatial location $t_0 \in \mathbb{R}^3$, the convolution feature output at this point is a stack of $[[\kappa * f](x_i)]_i$ with $x_i \in (S^{2'}, t_0) \subset X^\dagger$. As shown in Eq. (4), it involves feature gathering for \hat{f} at $s(x_i) \cdot x_{jk}$ for every $x_{jk} \in X' = S^{2'} \times \mathbb{R}^{3'}$. Denote an instance of $s(x_i) = (R_i, t_0)$, $x_{jk} = (R_j\mathbf{n}, t_k)$, then $s(x_i) \cdot x_{jk} = (R_iR_j\mathbf{n}, t_0 + R_it_k)$. If $\mathbb{R}^{3'}$ is closed under $SO(3)'$, then we have

$$\{R_it_k \mid t_k \in \mathbb{R}^{3'}\} = \mathbb{R}^{3'}, \forall R_i \in SO(3)', \quad (7)$$

which implies that the feature gathering for all $S^{2'}$ channels can be done once at the same set of spatial positions specified by $\mathbb{R}^{3'}$. An illustration is shown in Fig. 3. Without the symmetric kernel, the KPConv on group or quotient space looks like Fig. 3 (b), where the kernel is rotated for each $SO(3)'$ (for group KPConv) or $S^{2'}$ (for quotient KPConv) channel separately to gather input features. However, with symmetric kernels, as shown in Fig. 3 (c), the rotations keep the position of kernel points unchanged up to a permutation so that the feature-gathering step is simplified. Specifically, the number of spatial positions needed for feature gathering without symmetric kernels is $|SO(3)'||\mathbb{R}^{3'}|$ for group KPConv or $|S^{2'}||\mathbb{R}^{3'}|$ for quotient KPConv, while the number

Table 1. Efficiency comparison in terms of GPU memory consumption and the computation speed between EPN [2] and our method on three tasks. Two numbers are reported for *training/inference* respectively. ↓ means lower is better. ↑ means higher is better. The best is shown in bold font.

Tasks	ModelNet40 Pose		ModelNet40 Classification		3DMatch Keypoint Matching	
Methods	Memory (GB) ↓	Speed (fps) ↑	Memory (GB) ↓	Speed (fps) ↑	Memory (GB) ↓	Speed (fps) ↑
EPN [2]	22.2 / 16.9	1.1 / 1.6	13.4 / 12.7	1.9 / 1.5	37.4 / 8.5	0.6 / 3.1
<i>Ours</i> (w/o symmetric kernels)	4.8 / 3.7	5.1 / 10.1	4.1 / 3.2	7.8 / 7.8	7.5 / 2.8	2.6 / 16.7
<i>Ours</i> (w/ symmetric kernels)	4.3 / 2.8	6.7 / 11.1	3.9 / 2.7	9.1 / 10.3	6.5 / 2.4	3.7 / 23.6

is $|\mathbb{R}^{3'}|$ with symmetric kernels. The convolution kernels in Fig. 3 is an abstract illustration, while the actual symmetric kernel in our work is visualized in Fig. 4.

3.2.4 Choices of the discretization of $SO(3)$ and S^2

In Sec. 3.1.1, we mentioned that the discretization of $SO(3)$ is the rotation group that respects the symmetry of a Platonic solid. There are five types of Platonic solids, corresponding to 3 finite rotation groups, as shown in Fig. 5. Choosing $SO(3)'$ to be any of them is valid, representing a discretization of $SO(3)$ to different resolutions. The different Platonic solids with the same rotation group $SO(3)'$ represents different discretizations of $S^{2'}$ and $SO(2)'$.

If we use a small $SO(3)'$ (for example, T), the strategy of using $\mathbb{R}^{3'} = \{rS^{2'} \cup 0 \mid r > 0\}$ could be problematic because the number of kernel points could be too few to learn representative features. In this case, we are free to design the kernel points differently, as long as they are *closed* under $SO(3)'$. For example, one may add kernel points at the center of all edges and/or faces. One may even use a combination of several polyhedrons with different radii r .

In this paper, we only choose the icosahedron as the Platonic solid to conduct experiments for three reasons. First, it has the finest discretization of $SO(3)$. Second, it has a smaller size of $S^{2'}$ compared with the dodecahedron, maximizing the benefit of working with quotient features. Third, it is consistent with existing methods [1, 2], enabling direct comparison with the baselines.

3.2.5 Other aspects of the network

We use element-wise scalar nonlinearity (ReLU and leaky ReLU) in the network. The spatial pooling is done by subsampling the input points and aggregating the features of neighboring input points to the subsampled points. Batch normalization is applied to normalize over the batch, $S^{2'}$, and \mathbb{R}^3 dimensions. All these choices follow the common practice of conventional CNNs (KPCConv [35]) and group convolutions (EPN [2]). We also adopted the group attentive pooling in EPN to pool over the $S^{2'}$ dimension and generate $SO(3)'$ -invariant features.

Depending on the specific task in the experiment, the prediction head has a slightly different design, composing

the last few layers of the network. The loss functions are inherited from EPN [2], including cross-entropy loss for classifications, L2 loss for residual pose regression, and batch-hard triplet loss for keypoint matching. We refer to the appendix for more details about the prediction heads and the loss functions.

3.3. Relation to existing work

In the context of literature on group-equivariant neural networks, our method is under the theoretical framework of *equivariant CNNs on homogeneous spaces* [8, 24, 46]. Our work is a new form of realization when working with 3D point clouds and adapting the convolution structure of KPCConv. We explore a balance of simplicity, efficiency, and expressiveness by finding the proper quotient space and discretization. Our method is an extension of group convolutions, leveraging their clean structure enabled by discretization. Our method can also be viewed as a steerable CNN with homogeneous space $S^2 \times \mathbb{R}^3$ and stabilizer subgroup $SO(2)$ with scalar-type features or with homogeneous space \mathbb{R}^3 and stabilizer subgroup $SO(3)$ with S^2 features. The proposed work paves the way for efficiency improvement using quotient representation learning on finite groups. We also emphasize the group-variant side of equivariant models with the permutation layer to distinguish rotations, while existing work focuses on the benefit of getting group-invariant features from equivariant models.

4. Experiments

Our major baseline to compare with is EPN [2], a state-of-the-art group convolution network, as we are similar in several ways. Both are KPCConv-style convolutions with $SE(3)$ -equivariance. Both use the icosahedral rotation group \mathcal{I} to discretize $SO(3)$. However, EPN has features defined on $\mathcal{I} \times \mathbb{R}^3$, compared with $S^{2'} \times \mathbb{R}^3$ in our method.

Two datasets, ModelNet40 [44] and 3DMatch [47], are used in the experiments. ModelNet40 is composed of 3D CAD models of 40 categories of objects. 3DMatch is a real-scan dataset of indoor scenes. For the ModelNet40 dataset, we conduct the classification and pose estimation tasks. For the 3DMatch dataset, we conduct the keypoint matching task. These tasks are also studied in EPN [2].

In Tab. 1, we list the GPU memory consumption and run-

Table 2. Experimental result of object classification on ModelNet40. The best is **bolded**. The best in equivariant models is underlined. *Noisy*: test using input with random translation, scaling, jittering, and dropout. *Clean*: test without above processing. *SO(3)*: random rotations. *Id*: no rotation. *ico*: random rotations in \mathcal{T} . ESCNN works with voxelized data, thus not having *Noisy* results with point-wise augmentations. FLOPs are counted using `fvcore`, which does not support `DGL` used in TFN and SE(3)-T, thus left blank. The efficiency comparison uses the same batch size as in Tab. 1 (see appendix).

Type	Column #	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10	#11	#12	#13	#14
	Metrics	ModelNet40 classification performance metric: Acc (%) \uparrow										Efficiency metrics			
Non-equiv	Train rotation	SO(3)										Memory (GB) \downarrow train/test	Speed (fps) \uparrow train/test	FLOPs per training batch (G)	Trainable params (M)
	Test rotation	SO(3)				Id				ico					
	Test condition	Noisy	Clean	Noisy	Clean	Noisy	Clean	Noisy	Clean	Noisy	Clean				
Non-equiv	KPConv [35]	75.71	81.17	76.88	82.02	12.66	12.22	89.46	91.85	-	-	0.18/0.25	17.24/18.48	0.74	1.69
	PointNet++ [31]	81.31	84.88	83.03	85.76	13.44	13.76	90.67	91.55	-	-	1.03/0.5	7.06/7.43	10.59	1.48
	DGCNN [38]	79.86	84.77	82.54	85.62	15.36	17.26	91.00	92.18	-	-	1.96/1.69	17.90/17.50	32.65	1.81
	PT [48]	78.51	79.68	78.83	79.68	16.53	16.77	86.98	88.10	-	-	4.16/5.82	4.80/5.03	220.99	9.58
	CurveNet [45]	84.72	88.33	85.82	88.94	17.34	18.03	91.53	92.63	-	-	1.20/0.33	5.61/7.54	4.22	2.14
	PCT [17]	86.91	89.10	87.64	89.83	16.33	17.95	91.69	92.87	-	-	1.18/0.80	8.54/ 20.85	27.41	2.87
Equivariant	ESCNN [1]	-	82.40	-	77.73	-	29.03	-	88.94	-	-	4.71/5.50	3.79/8.38	428.61	0.46
	TFN [36]	58.27	62.64	58.06	62.64	57.50	62.28	59.20	62.28	-	-	15.95/8.7	1.83/6.08	-	0.06
	SE(3)-T [15]	60.37	66.29	61.18	66.29	44.61	50.53	44.81	50.53	-	-	19.39/9.65	1.54/5.13	-	0.11
	EPN [2]	84.63	87.84	85.34	88.83	30.99	32.32	90.03	91.61	90.04	91.60	13.40/12.72	1.86/1.49	58.95	3.06
	<i>Ours</i> (w/ GA pooling [2])	85.04	87.51	85.49	87.63	41.46	44.43	89.73	90.50	90.00	90.50	<u>3.95/2.70</u>	<u>9.21/10.37</u>	170.68	2.53
	<i>Ours</i> (w/ permutation)	86.99	88.62	88.21	89.62	39.54	42.78	<u>90.66</u>	<u>91.77</u>	90.68	91.77	<u>3.95/2.70</u>	9.09/10.28	170.77	2.65

ning speed of our method and EPN [2] in the three tasks. The comparisons are under the same input size, number of feature channels, and number of network layers. The separable convolutions on $SO(3)$ and \mathbb{R}^3 in EPN are together considered as one layer. The numbers are not comparable between training and inference because the batch size could be different (see the appendix). The specific configurations in each experiment are introduced later. All experiments are run on a single NVIDIA A40 GPU. Our network with quotient space convolutions is much smaller and runs much faster in all three tasks, indicating the potential value for real applications. We boost efficiency without sacrificing performance, as is shown later.

Ablation study: We list two rows of our results in Tab. 1 to show the effect of efficient feature gathering enabled by the symmetric kernels. The results of *without symmetric kernels* are generated using the same kernel points as *with symmetric kernels*, but ignoring the fact that they are symmetric to rotations and gathering features at $|S^{2'}||\mathbb{R}^{3'}|$ locations, instead of at $|\mathbb{R}^{3'}|$ locations and permuting them. The efficient feature gathering brings further efficiency improvements, especially in terms of computational speed.

4.1. Object classification on ModelNet40

For this task, given a point cloud of an object, the network predicts its category. The evaluation metric is classification accuracy (Acc). In this experiment, we show an extensive efficiency comparison with more existing networks, equivariant and non-equivariant. We also examine a wide combination of input conditions in training and testing to show the effect of equivariance and robustness against input imperfections. All models are trained with data augmentation, including random translation, scaling, jittering, and dropout.

Our method has outstanding performance as shown in Tab. 2. In columns #1-4, all methods are trained with rota-

Table 3. Pose estimation on ModelNet40. Mean, median, and max angular errors are calculated over the test set. Statistics (average and standard deviation) over 10 test runs are shown to account for the randomness.

Metrics	Mean ($^\circ$) \downarrow		Median ($^\circ$) \downarrow		Max ($^\circ$) \downarrow	
	Avg μ	SD σ	Avg μ	SD σ	Avg μ	SD σ
KPConv [35]	13.99	1.53	10.70	0.81	115.19	57.84
EPN [2]	1.10	0.20	1.36	0.13	7.06	2.52
<i>Ours</i>	1.20	0.08	0.96	0.05	6.71	1.28

tional augmentation. Our method performs the best among listed equivariant models and is on par with PCT [18] among all models. Our method is intended to work with rotational augmentations so that the equivariance gap caused by discretization can be interpolated through training. However, we still experiment with training without rotational augmentation, as shown in columns #5-10. From columns #9-10, we can verify the equivariance to the icosahedral rotation group of our model. Columns #5-6 show the effect of discretization on the continuous $SO(3)$ if no interpolation is trained, in which case the performance lands between non-equivariant models and continuously $SO(3)$ -equivariant models. The efficiency of our method outperforms all equivariant baselines and is similar to non-equivariant models. The FLOPs and number of trainable parameters are also listed for reference. We found that FLOPs do not strictly correlate with running speed, which could be due to different memory access costs and parallelism.

Ablation study: Since this task requires SE(3)-invariant features, there are two options in our network: either using the group-attentive pooling (GA pooling) introduced in EPN [2] to pool over the $S^{2'}$ dimensions or using the permutation layer to find the canonical permutation of $S^{2'}$ dimensions. Either way, the canonical pose of objects in ModelNet40 can be used for supervision. Tab. 2 shows that the permutation layer yields better performance with negligible computational overhead. The reason could be that the permutation layer as in Eq. (1) stacks features from $S^{2'}$

Table 4. Experiment result of keypoint matching on the 3DMatch dataset. The numbers are the average recall (%), and the higher, the better. Notation * represents the result with the given point normal information.

	SHOT [37]	3DM [47]	CGF [23]	PPFN [12]	PPFF [11]	3DSN [16]	Li [28]	Li [28]*	EPN [2]	Ours
Kitchen	74.3	58.3	60.3	89.7	78.7	97.5	92.1	99.4	99.0	99.4
Home 1	80.1	72.4	71.1	55.8	76.3	96.2	91.0	98.7	99.4	98.7
Home 2	70.7	61.5	56.7	59.1	61.5	93.2	85.6	94.7	96.2	96.6
Hotel 1	77.4	54.9	57.1	58.0	68.1	97.4	95.1	99.6	99.6	99.1
Hotel 2	72.1	48.1	53.8	57.7	71.2	92.8	91.3	100.0	97.1	98.1
Hotel 3	85.2	61.1	83.3	61.1	94.4	98.2	96.3	100.0	100.0	100.0
Study	64.0	51.7	37.7	53.4	62.0	95.0	91.8	95.5	96.2	95.2
MIT Lab	62.3	50.7	45.5	63.6	62.3	94.1	84.4	92.2	93.5	90.9
Average	73.3	57.3	58.2	62.3	71.8	95.6	91.0	97.5	97.6	97.3

and thus preserves the information better, compared with weighted averaging over the S^{2l} dimension as done in GA pooling.

4.2. Pose Estimation on ModelNet40

In this experiment, the network takes a pair of point clouds of an object and predicts the relative rotation between them. To avoid the pose ambiguity of objects with symmetric rotational shapes, only the airplane category is used in this experiment, with 626 models in the training set and 100 models in the test set. A point cloud is generated by randomly subsampling 1,024 points on the surface, and it is randomly rotated to form a pair.

The experimental result is shown in Tab. 3. We achieved similar rotation estimation accuracy to EPN [2] overall. While our mean error is slightly larger, the lower median error, max error, and standard deviations show that our method delivers a more reliable registration. It could imply that representing rotations as a permutation of features is more robust than representing them as a single element in the feature map. Besides, the equivariant networks outperform the non-equivariant KPConv [35] by a large margin.

4.3. Keypoint matching on 3DMatch

In this task, patches of point clouds extracted locally around keypoints in a large, dense scan are input to the network, and each is mapped to a feature vector of 64-dimension as the keypoint descriptor. Each patch has 1,024 points. Then we evaluate the average recall of keypoint correspondence across different scans through nearest neighbor search in the feature space, as proposed in PPFNet [12].

This experiment’s performance in Tab. 4 indicates the capability of learning distinctive and rotation-invariant features for local patches of point clouds. Though not achieving the best in the list, our method delivers comparable performance to the top methods using only a fraction of the computational resources as EPN [2] (see Tab. 1). We use the GA pooling layer [2] in this experiment because the permutation layer requires supervision on the pose, while

a canonical pose is not defined for local patches of point clouds, and GA pooling works with or without the pose supervision. However, the result shows that GA pooling over the $|S^{2l}|$ features also provides distinctive features for keypoint matching. This part may be further improved by taking the information of the global scan [12] or the matching scan [20] into consideration, in which case the permutation layer may get hints on the optimal permutation from the larger context. This topic goes beyond the focus of this paper and is left for future work.

5. Conclusion

This paper presents a new design of SE(3)-equivariant point cloud convolution network, which is efficient, simple, and expressive simultaneously by working with feature maps defined on the quotient space $S^2 \times \mathbb{R}^3$ associated with the stabilizer SO(2). We further improve the efficiency of the convolutions by designing the kernel points to be symmetric to the discretized rotation group $SO(3)'$. Moreover, we propose a permutation layer to recover $SO(3)'$ information from S^{2l} dimensions of the features so that the network can detect SO(3) rotations. Experiments show that our network delivers state-of-the-art performance in multiple tasks while consuming only a fraction of memory and computation resources as a group-convolution network with similar performance. Our method can open exciting opportunities to introduce the SE(3)-equivariance property to mainstream point cloud networks for various tasks.

This work also has limitations. We do not outperform EPN [2] in the keypoint matching task, implying that the network, especially the permutation layer, needs improvement when dealing with inputs without a clear pose definition. Other possibilities for the network design also remain open. For example: What if we use a non-scalar type of feature on the quotient space? How to further alleviate the impact of the discretization of a continuous group? From a general perspective, extending the discretized quotient space convolution strategy to other groups is also an attractive direction for future work.

References

- [1] Gabriele Cesa, Leon Lang, and Maurice Weiler. A program to build E(N)-equivariant steerable CNNs. In *International Conference on Learning Representations*, 2022.
- [2] Haiwei Chen, Shichen Liu, Weikai Chen, Hao Li, and Randall Hill. Equivariant point network for 3D point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 14514–14523, 2021.
- [3] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
- [4] Taco Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Convolutional networks for spherical signals. *arXiv preprint arXiv:1709.04893*, 2017.
- [5] Taco Cohen, Maurice Weiler, Berkay Kicanaoglu, and Max Welling. Gauge equivariant convolutional networks and the icosahedral CNN. In *International conference on Machine learning*, pages 1321–1330. PMLR, 2019.
- [6] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR, 2016.
- [7] Taco S Cohen, Mario Geiger, and Maurice Weiler. Intertwiners between induced representations (with applications to the theory of equivariant neural networks). *arXiv preprint arXiv:1803.10743*, 2018.
- [8] Taco S Cohen, Mario Geiger, and Maurice Weiler. A general theory of equivariant CNNs on homogeneous spaces. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 32, 2019.
- [9] Taco S Cohen and Max Welling. Steerable cnns. *arXiv preprint arXiv:1612.08498*, 2016.
- [10] Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulénard, Andrea Tagliasacchi, and Leonidas J Guibas. Vector Neurons: A general framework for SO(3)-equivariant networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 12200–12209, 2021.
- [11] Haowen Deng, Tolga Birdal, and Slobodan Ilic. PPF-FoldNet: Unsupervised learning of rotation invariant 3D local descriptors. In *Proceedings of the European Conference on Computer Vision*, pages 602–618, 2018.
- [12] Haowen Deng, Tolga Birdal, and Slobodan Ilic. PPFNet: Global context aware local features for robust 3D point matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 195–205, 2018.
- [13] Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing convolutional neural networks for equivariance to lie groups on arbitrary continuous data. In *International Conference on Machine Learning*, pages 3165–3176. PMLR, 2020.
- [14] Marc Finzi, Max Welling, and Andrew Gordon Wilson. A practical method for constructing equivariant multilayer perceptrons for arbitrary matrix groups. In *International Conference on Machine Learning*, pages 3318–3328. PMLR, 2021.
- [15] Fabian Fuchs, Daniel Worrall, Volker Fischer, and Max Welling. SE(3)-transformers: 3D roto-translation equivariant attention networks. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 33:1970–1981, 2020.
- [16] Zan Gojcic, Caifa Zhou, Jan D Wegner, and Andreas Wieser. The perfect match: 3D point cloud matching with smoothed densities. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5545–5554, 2019.
- [17] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7:187–199, 2021.
- [18] Meng-Hao Guo et al. Pct: Point cloud transformer. *Computational Visual Media*, 2021.
- [19] Emiel Hooeboom, Jorn WT Peters, Taco S Cohen, and Max Welling. HexaConv. In *International Conference on Learning Representations*, 2018.
- [20] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. Predator: Registration of 3D point clouds with low overlap. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4267–4276, 2021.
- [21] Michael J Hutchinson, Charline Le Lan, Sheheryar Zaidi, Emilien Dupont, Yee Whye Teh, and Hyunjik Kim. LieTransformer: Equivariant self-attention for lie groups. In *International Conference on Machine Learning*, pages 4533–4543. PMLR, 2021.
- [22] Nicolas Keriven and Gabriel Peyré. Universal invariant and equivariant graph neural networks. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 32, 2019.
- [23] Marc Khoury, Qian-Yi Zhou, and Vladlen Koltun. Learning compact geometric features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 153–161, 2017.
- [24] Risi Kondor and Shubhendu Trivedi. On the generalization of equivariance and convolution in neural networks to the action of compact groups. In *International Conference on Machine Learning*, pages 2747–2755. PMLR, 2018.
- [25] Abhinav Kumar, Garrick Brazil, Enrique Corona, Armin Parchami, and Xiaoming Liu. Deviant: Depth equivariant network for monocular 3d object detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 664–683. Springer, 2022.
- [26] Leon Lang and Maurice Weiler. A Wigner-Eckart theorem for group equivariant convolution kernels. *arXiv preprint arXiv:2010.10952*, 2020.
- [27] Jongmin Lee, Byungjin Kim, and Minsu Cho. Self-supervised equivariant learning for oriented keypoint detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4847–4857, June 2022.
- [28] Lei Li, Siyu Zhu, Hongbo Fu, Ping Tan, and Chiew-Lan Tai. End-to-end learning local multi-view descriptors for 3D point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1919–1928, 2020.

- [29] Chien Erh Lin, Jingwei Song, Ray Zhang, Minghan Zhu, and Maani Ghaffari. Epn-netvlad: Se (3)-invariant place recognition for 3d point clouds. In *6th Annual Conference on Robot Learning*.
- [30] Lachlan E MacDonald, Sameera Ramasinghe, and Simon Lucey. Enabling equivariance for arbitrary Lie groups. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8183–8192, 2022.
- [31] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 30, 2017.
- [32] Nimrod Segol and Yaron Lipman. On universal equivariant set networks. *arXiv preprint arXiv:1910.02421*, 2019.
- [33] Ahyun Seo, Byungjin Kim, Suha Kwak, and Minsu Cho. Reflection and rotation symmetry detection via equivariant learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9539–9548, June 2022.
- [34] Anthony Simeonov, Yilun Du, Andrea Tagliasacchi, Joshua B Tenenbaum, Alberto Rodriguez, Pulkit Agrawal, and Vincent Sitzmann. Neural Descriptor Fields: SE(3)-equivariant object representations for manipulation. In *Proceedings of the IEEE International Conference on Robotics and Automation*, pages 6394–6400. IEEE, 2022.
- [35] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [36] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- [37] Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique shape context for 3d data description. In *Proceedings of the ACM workshop on 3D object retrieval*, pages 57–62, 2010.
- [38] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (TOG)*, 38(5):1–12, 2019.
- [39] Maurice Weiler and Gabriele Cesa. General E(2)-equivariant steerable CNNs. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 32, 2019.
- [40] Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco S Cohen. 3D steerable CNNs: Learning rotationally equivariant features in volumetric data. *Proceedings of the Advances in Neural Information Processing Systems Conference*, 31, 2018.
- [41] Marysia Winkels and Taco S Cohen. 3D G-CNNs for pulmonary nodule detection. *arXiv preprint arXiv:1804.04656*, 2018.
- [42] Daniel Worrall and Gabriel Brostow. CubeNet: Equivariance to 3D rotation and translation. In *Proceedings of the European Conference on Computer Vision*, pages 567–584, 2018.
- [43] Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhambetov, and Gabriel J Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5028–5037, 2017.
- [44] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1912–1920, 2015.
- [45] Tiange Xiang, Chaoyi Zhang, Yang Song, Jianhui Yu, and Weidong Cai. Walk in the cloud: Learning curves for point clouds shape analysis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 915–924, October 2021.
- [46] Yinshuang Xu, Jiahui Lei, Edgar Dobriban, and Kostas Daniilidis. Unified Fourier-based kernel and nonlinearity design for equivariant networks on homogeneous spaces. In *International Conference on Machine Learning*, pages 24596–24614. PMLR, 2022.
- [47] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3DMatch: Learning local geometric descriptors from RGB-D reconstructions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1802–1811, 2017.
- [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] Minghan Zhu, Maani Ghaffari, and Huei Peng. Correspondence-free point cloud registration with so (3)-equivariant implicit shape representations. In *Conference on Robot Learning*, pages 1412–1422. PMLR, 2022.

E2PN: Efficient SE(3)-Equivariant Point Network (Appendix)

Minghan Zhu
University of Michigan
minghanz@umich.edu

Maani Ghaffari
University of Michigan
maanigj@umich.edu

William A Clark
Cornell University
wac76@cornell.edu

Huei Peng
University of Michigan
hpeng@umich.edu

1. Preliminaries and notation

We first review some basic concepts in group theory and representation theory briefly. They are highly relevant for understanding the big picture from a theoretical perspective of our work and general equivariant deep learning literature.

Groups: A group G is a set equipped with a binary operation \cdot , satisfying the following conditions: (1) the set is closed under the operation: $x \cdot y \in G, \forall x, y \in G$; (2) the operation is associative: $x \cdot (y \cdot z) = (x \cdot y) \cdot z, \forall x, y, z \in G$; (3) there is an identity element e in the set such that $x \cdot e = e \cdot x = x, \forall x \in G$; (4) there is an inverse x^{-1} for each element x in the set such that $x \cdot x^{-1} = x^{-1} \cdot x = e$. For example, the integer set \mathbb{Z} is a group under the addition operator with identity 0 and inverse $-x$ for any $x \in \mathbb{Z}$. Sometimes we omit the \cdot notation.

Group actions and representations: We say a group G acts on a set X if any element g in G corresponds to a transformation $\rho(g)$ on X , i.e., $[\rho(g)](x) \in X, \forall x \in X$, such that $\rho(g_1) \circ \rho(g_2) = \rho(g_1 g_2), \forall g_1, g_2 \in G$, where \circ denotes function compositions, and $[\rho(e)](x) = x, \forall x \in X$. When X is a linear space and $\rho(g)$ is linear, we say ρ is a (linear) *representation* of G in X . When X is n (finite)-dimensional linear space, we have a representation $\rho: G \rightarrow \text{GL}(n)$, i.e., we can write $[\rho(g)](x)$ as $\rho(g)x$, where $\rho(g)$ takes the form of n -by- n invertible matrices and acts by multiplication on the left. For example, the representations of 2D rotations $\text{SO}(2)$ in \mathbb{R}^2 are the 2-by-2 orthonormal matrices with determinant 1. We also use (ρ, X) as a shorthand to denote the representation and space on which it acts.

Equivariance: Given spaces V_1 with representation ρ_1 of G and V_2 with representation ρ_2 of G , we say a mapping $\phi: V_1 \rightarrow V_2$ is G -equivariant if $\phi \circ \rho_1(g) = \rho_2(g) \circ \phi, \forall g \in G$. A G -equivariant linear map is also called an *intertwiner*. The space of intertwiners is denoted $\text{Hom}_G(\rho_1, \rho_2)$, homomorphisms of group representations ρ_1, ρ_2 of G .

Subgroups, cosets, and quotient spaces: A subgroup H of G is a subset of G that is also a group, denoted $H \leq G$.

For example, $\text{SO}(2) \leq \text{SO}(3)$. Given $H \leq G$ and $g \in G$, we can define a (left-)coset as $gH = \{gh|h \in H\}$. For a given H , all cosets are either equal or disjoint. Each coset is of the same size (contains the same number of elements), and they partition the whole group. The set of cosets forms a coset space (or *quotient space*) $G/H = \{gH|g \in G\}$. In short, a coset is both a subset in the group and an element in the quotient space.

Stabilizer subgroup: If a group G acts on set X by ρ , for $x \in X$, the stabilizer subgroup is defined as $\text{Stab}_G(x) \triangleq \{g \in G|\rho(g)x = x\}$. By definition, the stabilizer subgroup $\text{Stab}_G(e_{G/H})$ for the quotient space G/H is H .

Homogeneous spaces: Assume that a group G acts on a space X through action ρ , we call X a *homogeneous space* of G if G acts *transitively* on X , i. e., any two elements in X are connected by a group action, $\forall x_1, x_2 \in X, \exists g \in G$, s.t. $x_1 = \rho(g)x_2$. A quotient space G/H is a homogeneous space of G .

Induced representations: Here is an important known result [1, 4, 5]: given a representation ρ of subgroup H on vector space V , one can *induce* a representation $\pi = \text{Ind}_H^G \rho$ of G for the space of functions $\mathcal{F} = \{f: G/H \rightarrow V\}$. It provides a way to define group actions in function spaces, a foundation of the research on equivariant feature learning.

2. Definition of the section functions

In Sec. 3.2.2, we define the convolution in a homogeneous space as Eq. (4), using the section function $s: X \rightarrow G$, mapping an element in the quotient space to a group element in the corresponding coset, i.e.,

$$s(x)H = x, \forall x \in X = G/H \quad (8)$$

In our work, for the continuous case, $H = \text{SO}(2), G = \text{SO}(3), X = S^2$, and Eq. (8) can be rewritten as

$$s(x)\mathbf{n} = x, \forall x \in S^2 \quad (9)$$

Since there are generally multiple group elements in a coset, section functions are not unique. Thus we need to define the section function to make the convolution well-defined. For any $R \in \text{SO}(3)$, we can write $R = R_z(\alpha)R_y(\beta)R_z(\gamma)$ using Euler angles $\alpha \in [0, 2\pi), \beta \in [0, \pi], \gamma \in [0, 2\pi)$, and the coset it belongs to is $R\mathbf{n} = \{RR_z(\theta)|\theta \in [0, 2\pi)\} = \{R_z(\alpha)R_y(\beta)R_z(\gamma + \theta)|\theta \in [0, 2\pi)\} = R_z(\alpha)R_y(\beta)\mathbf{n}$. Thus a natural section from S^2 to $\text{SO}(3)$ is

$$s(R\mathbf{n}) \triangleq R_z(\alpha)R_y(\beta) \quad (10)$$

which removes the last z -axis rotation in z - y - z Euler angle rotations. In the discretized setup, Eq. (10) does not work because for $R\mathbf{n} \in S^{2'} \subset S^2$, $s(R\mathbf{n}) \in \text{SO}(3)$ may not be in $\text{SO}(3)'$. In this case, we just arbitrarily select an element in each coset as the section so that $s' : S^{2'} \rightarrow \text{SO}(3)'$ satisfies Eq. (8). While the selection is arbitrary, it should be fixed once selected so that the behavior of the function is deterministic and consistent for different inputs.

3. The derivation of our proposed convolution

3.1. The equivariance of our convolution

The derivation of the convolution in this paper is mostly built upon [4]. We do not discover new theorems. Our result is an application of the existing theoretical results in a specification that is not previously discussed in the literature. Here we start from the conclusions in [4] and show how it leads to our design of convolutions.

We first need to introduce another concept. For any $g \in G, x \in G/H$,

$$(gs(x))H = g(s(x)H) = g(x) = (s(gx))H \quad (11)$$

meaning that $gs(x)$ and $s(gx)$ are in the same coset, but $gs(x)$ and $s(gx)$ are not necessarily equal. We can relate these two using a function $h : G/H \times G \rightarrow H$ as:

$$gs(x) = s(gx)h(x, g) \quad (12)$$

i.e., $h(x, g) \triangleq s(gx)^{-1}gs(x)$. This function h describes how the representative group element twists beyond jumping to another coset when applied with another group element, therefore heavily relying on s . We may denote it as h_s , but we will go with h in the following since we already selected and fixed s in Sec. 2.

With this h function, we can write down the form of induced representations. Given a space of functions $\mathcal{F} = \{f : G/H \rightarrow V\}$, assuming $\rho : H \rightarrow \text{GL}(V)$ a representation of subgroup H in V , we define $\pi = \text{Ind}_H^G \rho : G \rightarrow \text{GL}(\mathcal{F})$ as:

$$[\pi(g)f](x) \triangleq \rho(h(g^{-1}x, g))f(g^{-1}x) \quad (13)$$

It is shown in [4] that Eq. (13) is a valid representation. Denote $\mathcal{F}_1 = \{G/H \rightarrow V_1\}$ and $\mathcal{F}_2 = \{G/H \rightarrow V_2\}$ with

representations (ρ_1, V_1) and (ρ_2, V_2) on H , any linear mapping $\mathcal{F}_1 \rightarrow \mathcal{F}_2$ equivariant to the induced representations $\text{Ind}_H^G \rho_1$ and $\text{Ind}_H^G \rho_2$ can be written as a cross-correlation with a twist:

$$[\kappa * f](x) = \int_{G/H} \kappa(s(x)^{-1}y)\rho_1(h(y, s(x)^{-1}))f(y)dy \quad (14)$$

where the space \mathcal{K}_C of valid kernels $\kappa : G/H \rightarrow \text{Hom}(V_1, V_2)$ is equivalent to the space:

$$\begin{aligned} \mathcal{K}_D &= \{\bar{\kappa} : H \backslash G/H \rightarrow \text{Hom}(V_1, V_2) | \\ \bar{\kappa}(x) &= \rho_2(h)\bar{\kappa}(x)\rho_1^x(h)^{-1}, \forall x \in H \backslash G/H, h \in H^{\eta(x)H}\} \end{aligned} \quad (15)$$

where $H \backslash G/H = \{HgH | g \in G\}$ is the set of double cosets in which $HgH = \{h_1gh_2 | h_1, h_2 \in H\}$ is called a double coset. $\eta : H \backslash G/H \rightarrow G$ is a section function for double cosets. $H^{\eta(x)H} = \{h \in H | h\eta(x)H = \eta(x)H\}$ is the set of stabilizers for double coset $x \in H \backslash G/H$, and ρ_1^x is the representation of $H^{\eta(x)H}$ defined as $\rho_1^x(h) = \rho_1(\eta(x)^{-1}h\eta(x))$ for $h \in H^{\eta(x)H}$.

While the above looks a bit involved, recall that in this work, we use scalar-type features, meaning that we choose the trivial identity representation for the subgroup $H = \text{SO}(2)$, i.e., $\rho_1(h) = \text{Id}_{V_1}, \rho_2(h) = \text{Id}_{V_2} \forall h \in H$, which simplifies the equations. The induced representation $\pi = \text{Ind}_H^G \rho : G \rightarrow \text{GL}(\mathcal{F})$ is in the form:

$$[\pi(g)f](x) = f(g^{-1}x), \forall g \in G, x \in G/H \quad (16)$$

The convolution in Eq. (14) now looks like:

$$\begin{aligned} [\kappa * f](x) &= \int_{G/H} \kappa(s(x)^{-1}y)f(y)dy \\ &= \int_{G/H} \kappa(y)f(s(x)y)dy \end{aligned} \quad (17)$$

which is consistent with Eq. (4) in the main paper. The equivalent space of kernels is:

$$\mathcal{K}_D = \{\bar{\kappa} : H \backslash G/H \rightarrow \text{Hom}(V_1, V_2)\} \quad (18)$$

3.2. The specific form of our kernel

In this paper, we work with $G = \text{SE}(3)$ and $H = \text{SO}(2)$. In the following, we derive G/H and $H \backslash G/H$ in this setup since they do not appear commonly in the literature.

The group $\text{SE}(3) = \text{SO}(3) \ltimes \mathbb{R}^3$ is the semi-direct product of $\text{SO}(3)$ and \mathbb{R}^3 (the latter is a normal subgroup). We can denote a group element of $\text{SE}(3)$ as (R, t) where $R \in \text{SO}(3), t \in \mathbb{R}^3$, such that the group action \cdot is defined as $(R_1, t_1) \cdot (R_2, t_2) = (R_1R_2, R_1t_2 + t_1)$, and accordingly the group inverse is defined as $(R, t)^{-1} = (R^{-1}, -R^{-1}t)$.

The group $\text{SO}(3)$ can be written as a subgroup of $\text{SE}(3)$ as $(R, 0)$. Using Euler angles we have $\forall R \in$

$\text{SO}(3), \exists \alpha \in [0, 2\pi), \beta \in [0, \pi], \gamma \in [0, 2\pi)$, such that $R = R_z(\alpha)R_y(\beta)R_z(\gamma)$, where R_z represents rotation around the z -axis, and similarly for R_y . We also have $\text{SO}(2) \cong \{(R_z(\gamma), 0) \in \text{SE}(3) | \gamma \in [0, 2\pi)\}$.

Therefore, a left coset of $H = \text{SO}(2)$ in $G = \text{SE}(3)$ is $gH = \{gh | h \in H\} = \{(R_g, t_g) \cdot (R_z(\gamma_h), 0) | \gamma_h \in [0, 2\pi)\} = \{(R_z(\alpha_g)R_y(\beta_g)R_z(\gamma_g + \gamma_h), t_g) | \gamma_h \in [0, 2\pi)\}$, meaning that a left coset can be parameterized by α_g, β_g, t_g . Then the set of left cosets G/H is homeomorphic to the Cartesian product $S^2 \times \mathbb{R}^3 = \{(R_z(\alpha)R_y(\beta)\mathbf{n}, t) | \alpha \in [0, 2\pi), \beta \in [0, \pi], t \in \mathbb{R}^3\}$, where S^2 is the surface of a sphere, $\mathbf{n} = t(0, 0, 1)$ is the unit vector pointing to the north pole. Here we abuse the notation $(x\mathbf{n}, y)$ as an ordered pair in the set $S^2 \times \mathbb{R}^3$. It can be understood as a point x on a sphere centered at some point y in \mathbb{R}^3 . The group $G = \text{SE}(3)$ acts on G/H by left multiplication: $(R_g, t_g)(R\mathbf{n}, t) = (R_gR\mathbf{n}, R_g t + t_g)$.

We further investigate the double coset space $H \backslash G / H$. An element in the set can be written as $HgH = \{h_1gh_2 | h_1, h_2 \in H\} = \{(R_z(\alpha_g + \gamma_{h_1})R_y(\beta_g)R_z(\gamma_g + \gamma_{h_2}), R_z(\gamma_{h_1})t_g) | \gamma_{h_1}, \gamma_{h_2} \in [0, 2\pi)\}$. We can use $t(x, y, z)$ to specify the coordinate of an element in \mathbb{R}^3 , and we can always rewrite $t(x, y, z) = R_z(\gamma_t)t(r_g, 0, z_g)$ where $r_g = \sqrt{x^2 + y^2} \geq 0$ and $\gamma_t = \arctan 2(y, x)$. Then we can rewrite

$$HgH = \{(R_z(\alpha_g + \gamma_{h_1})R_y(\beta_g)R_z(\gamma_g + \gamma_{h_2}), R_z(\gamma_{h_1} + \gamma_t)t(r_g, 0, z_g)) | \gamma_{h_1}, \gamma_{h_2} \in [0, 2\pi)\} \quad (19)$$

Let us rename $\gamma_1 \triangleq \gamma_{h_1} + \gamma_t, \theta_g \triangleq \alpha_g - \gamma_t, \gamma_2 \triangleq \gamma_g + \gamma_{h_2}$, then we have

$$HgH = \{(R_z(\theta_g + \gamma_1)R_y(\beta_g)R_z(\gamma_2), R_z(\gamma_1)t(r_g, 0, z_g)) | \gamma_1, \gamma_2 \in [0, 2\pi)\} \quad (20)$$

Now it is clear that an element in $H \backslash G / H$ can be determined by four parameters $(\theta_g, \beta_g, r_g, z_g)$, with $\theta_g \in [0, 2\pi), \beta_g \in [0, \pi], r_g \geq 0, z_g \in \mathbb{R}$. We denote an element in $H \backslash G / H$ as $HgH(\theta_g, \beta_g, r_g, z_g)$. We have $H \backslash G / H \cong S^2 \times \mathbb{R}^+ \times \mathbb{R}$. Geometrically, each point (r_g, z_g) in the $\mathbb{R}^+ \times \mathbb{R}$ plane corresponds to a circle around the z -axis with radius r_g at height z_g . On the other hand, θ_g, β_g parameterizes a point on the sphere S^2 .

It follows that we can define a kernel $\bar{\kappa} \in \mathcal{K}_{\mathcal{D}} = \{\bar{\kappa} : S^2 \times \mathbb{R}^+ \times \mathbb{R} \rightarrow \text{Hom}(V_1, V_2)\}$, and then injectively map it to $\kappa \in \mathcal{K}_{C_0} = \{\kappa : S^2 \times \mathbb{R}^3 \rightarrow \text{Hom}(V_1, V_2)\}$. The fact that $S^2 \times \mathbb{R}^+ \times \mathbb{R} \subsetneq S^2 \times \mathbb{R}^3$ implies that $\mathcal{K}_{\mathcal{D}} \cong \mathcal{K}_{\mathcal{C}} \subsetneq \mathcal{K}_{C_0}$. In other words, there is a certain constraint on \mathcal{K}_{C_0} to form the actual set of valid equivariant kernels $\mathcal{K}_{\mathcal{C}}$. As shown in [4], the general form of the constraint can be written as:

$$\kappa(hx) = \rho_2(h)\kappa(x)\rho_1(h(x, h)^{-1}), \quad (21)$$

for $\kappa \in \mathcal{K}_{\mathcal{C}}, x \in G/H, h \in H$. As discussed in Sec. 3.1, ρ_1 and ρ_2 are both identity; thus Eq. (21) becomes $\kappa(hx) = \kappa(x)$, which is equivalent to Eq. (6) in the main paper in our specific case.

4. Equivariance of element-wise non-linear layers and normalization layers

For a feature map of shape $B \times C \times N \times A$, where B is the batch size, C is the feature channel, N is the number of spatial points in \mathbb{R}^3 , A corresponds to the spherical coordinates in $S^{2'}$, a batch normalization (BatchNorm [6]) is to calculate the mean and variance across the $B \times N \times A$ channels and apply a constant scaling factor and shift for each $B \times N \times A$ tensor. For instance normalization (InstanceNorm [8]), one just need to change $B \times N \times A$ to $N \times A$. In either case, consider the feature map as a function $f : N \times A \rightarrow \mathbb{R}^{B \times C}$, a BatchNorm or InstanceNorm (denoted \mathcal{N}) behaves like an element-wise operation, i.e.,

$$[\mathcal{N} \cdot f](x) = af(x) + b = \mathcal{N} \cdot f(x), \quad \forall x \in S^{2'} \times \mathbb{R}^3 \quad (22)$$

since a, b are constant vectors. Here \cdot denotes applying a transformation.

Recall that our induced representation Eq. (16) is in a similar form of a regular representation, which is realized by a change of coordinate without modifying the function value. Such a representation π is commutative with element-wise operations (denoted as \mathcal{E}):

$$\begin{aligned} [(\pi(g) \circ \mathcal{E}) \cdot f](x) &= [\pi(g) \cdot (\mathcal{E} \cdot f)](x) = [\mathcal{E} \cdot f](g^{-1}x) = \mathcal{E} \cdot f(g^{-1}x) \\ &= \mathcal{E} \cdot [\pi(g) \cdot f](x) = [\mathcal{E} \cdot (\pi(g) \cdot f)](x) \\ &= [(\mathcal{E} \circ \pi(g)) \cdot f](x), \forall g \in G \end{aligned} \quad (23)$$

Or we can say:

$$\pi(g) \circ \mathcal{E} = \mathcal{E} \circ \pi(g), \forall g \in G \quad (24)$$

It shows that element-wise non-linear layers like ReLU and normalization layers, including BatchNorm and InstanceNorm, are G -equivariant.

5. Prediction heads and loss functions

5.1. Pose estimation task

The pose (rotation) estimation task is fulfilled with a prediction head designed as shown in Fig. 1. The inputs to the prediction head for each pair of point clouds are two $S^{2'} \times C$ features, where C is the number of feature channels. We call the $S^{2'}$ coordinates *anchors* in this section. We apply $R_i \in \mathcal{I}$ to the second point-cloud mentally, corresponding to 60 permutations of the anchors for f_2 . If the two point clouds are different exactly by a rotation in \mathcal{I} , then

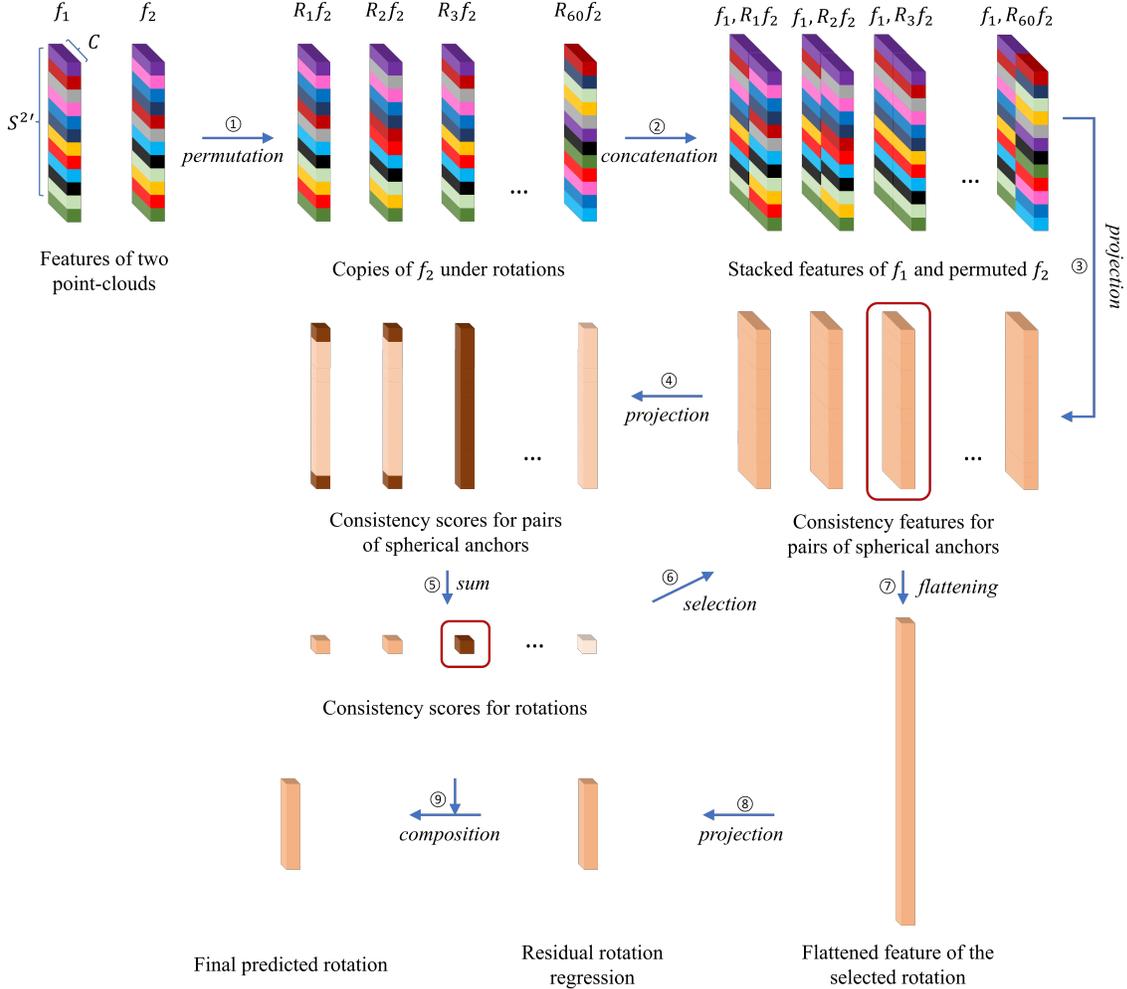


Figure 1. Illustration of the prediction head for rotation estimation. The numbers show the sequence of operations. The colors in the top row correspond to different spherical anchors. The shade of color after step 4 and after step 5 represents the matching scores for pairs of anchors. Darker means higher.

one of the permuted f_2 should be exactly the same as f_1 . We stack the features together and use several linear layers to find the match. Notice that the matching is defined as a binary classification problem for each pair of anchors instead of a multi-class classification problem for the overall feature corresponding to a certain rotation. It aligns better with the underlying geometry because a subset of anchors may align even under a wrong rotation since any rotations in $SO(2)$ keep the north-pole and south-pole vertices static. We can find the correct rotation class by summing over all anchor pairs and picking the rotation with the highest overall matching score. After finding the correct permutation (equivalent to the $R_i \in \mathcal{I} \subset SO(3)$), we flatten the feature and regress the residual rotation using quaternions in a way similar to [3].

Accordingly, the loss functions are the binary cross entropy loss for anchor-pair matching and L2 loss on the residual rotation regression.

5.2. Object classification task

For the classification task, we follow a similar philosophy as the rotation estimation task. Here we do not have a pair of inputs from which to find the relative rotation. Therefore, we imagine that there is a *reference* object for each category, with a canonical permutation of the features representing the underlying canonical pose. We learn the features of the reference object in each class and use them to classify input point clouds.

The core learnable parameter is the reference feature X of shape $S^{2'} \times C \times N$ where N is the number of object classes. We can denote $X_n \in \mathbb{R}^{S^{2'} \times C}$ as the reference feature for object class n . We directly calculate the inner product between the reference features and the permuted input features. The score of rotations under every object-class hypothesis is calculated by summing over the inner products across all anchors. The score of each class is de-

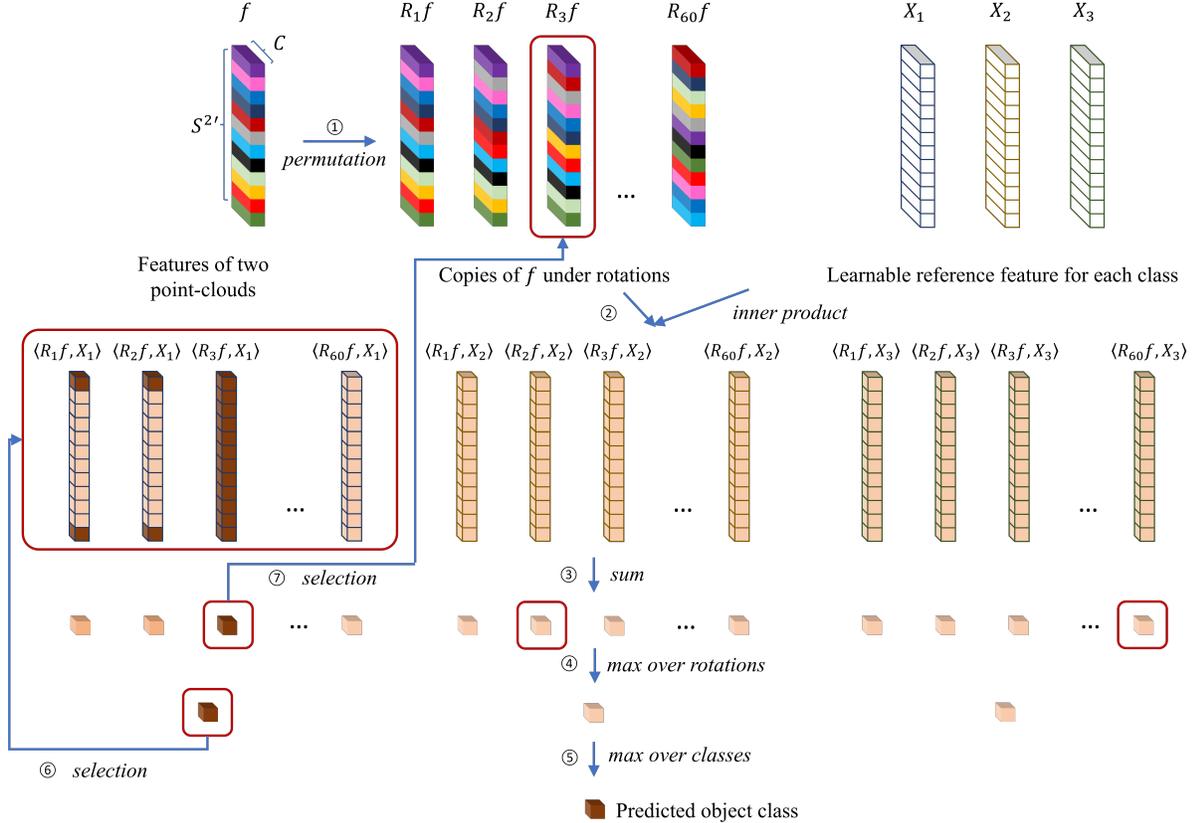


Figure 2. Illustration of the prediction head for object classification. The numbers show the sequence of operations. The solid colors in the top row correspond to different spherical anchors. The line colors on the top right represent different semantic classes. Here we only use three classes for illustration. The shade of color after step 2 represents the scores. Darker means higher. $\langle \cdot, \cdot \rangle$ denotes the inner product of finite-dim vectors.

defined as the maximum rotation score in each class hypothesis. We first generate the final prediction of the object class using the class score. Then we go back to the inner product matrix corresponding to this class and use it as the anchor-matching score prediction. Finally, the best rotation in this class is used to retrieve the specific permutation of the input feature, which forms a rotation-invariant feature for this object.

The basic assumption here is that only the correct object class yields a high-quality matching under the actual rotation. Thus we first solve for the classification and then determine the optimal rotation only in this class, which is used for generating rotation-invariant features.

The loss functions applied are the cross entropy loss for object classification and the binary cross entropy loss for anchor-matching prediction.

5.3. Keypoint matching task

For the keypoint matching task, there is no definition of a canonical pose for a local patch of points around a keypoint. Because the feature learning in this task does not involve the corresponding patch in another point cloud, we cannot

define relative poses as well. Therefore, we do not apply the permutation layer in this task. Instead, we simply follow the same design as in [3] using GA pooling, except that our attentive pooling is not over \mathcal{I} , but S^{2l} .

The loss function applied here is the batch-hard triplet loss, also consistent with [3].

6. More specifications in the experiments

The batch size used in the efficiency comparison in Tab. 1 is specified in Tab. 5. For the keypoint matching task, the number of global scans processed (n_g) and the number of local patches extracted from each global scan (n_l) define the input size. We use $n_g \times n_l$ as the notation in Tab. 5.

The training optimizer and learning rate schedule follow the default setup of EPN [3]. The number of feature channels (i.e., width) and the number of network layers (i.e., depth) also follow the settings in EPN, except that for the object classification task on ModelNet40, we reduced the backbone width by half compared with the original EPN setting (first layer width changed from 64 to 32, later layers in the backbone changed accordingly). The network width is the same across KPConv [7], EPN, and our E2PN in

Table 5. The batch sizes used in the efficiency comparison in terms of the GPU memory consumption and the running speed between EPN [3] and our method on three tasks as in Tab. 1.

Tasks	ModelNet40 Pose		ModelNet40 Classification		3DMatch Keypoint Matching	
Modes	Training	Inference	Training	Inference	Training	Inference
Batch size	8	8	12	24	1×16	8×24

our experiments, therefore their comparison remains valid. For all the three networks, the classification accuracy is not harmed by the width reduction.

Due to the page limit in the main paper, we only mentioned "all results are trained and tested with rotational augmentation" in the classification task section. Here we provide more details. During training, random rotations are applied, following the practice in EPN. EPN and our network are equivariant to a discretization of $SO(3)$; therefore, we apply rotation augmentations to let the network interpolate among the discretizations well and regress the residual rotation adding to the discretized rotations. During testing, we use fixed rotations for each test input so that results are repeatable and comparable.

We also provide more details here on the steerable CNN baseline based on ESCNN [2] in the classification task experiment. [2] established a general framework for steerable CNNs equivariant to $O(3)$ and its subgroups. It is relevant to the discussion in our paper and reported results on shape classification on the ModelNet10 dataset, which is a similar task to our experiments on ModelNet 40 dataset. Therefore we include it as one of our baselines. While the library for this work is open-sourced, the specific implementation of the network for the tasks mentioned in [2] is not provided. Therefore we implemented the network using the library according to the specifics stated in [2]. We also implemented the data conversion to transform the shapes into voxel grids as described in [2]. The Gaussian kernel radius in voxel generation and the learning rate are not specified in [2], so we did a hand tuning and reported the result under the best settings. We applied the final values $\sigma = 0.03$ for the Gaussian kernel in voxel generation and $lr = 10^{-3}$ as the initial learning rate. The multiplicity of irreps in the backbone is not specified either, and we found that using band-limited regular representations yields better results than using equal proportions of irreps. Thus we use multiple and-limited regular representations and align the total number of channels to the numbers specified in [2]. We implemented the $SO(3)$ equivariant version with frequency up to 3. Among the invariant maps discussed in [2], group pooling and norm pooling are implemented in their library. Group pooling is not recommended for continuous groups in [2]; therefore, we use norm pooling as the invariant layer.

The result is shown in Table 3 in the main paper. We did not compare the efficiency because the forms of input

(voxels vs. point clouds) and network structures (number of layers, channels, and connections) are both quite different. We can see that the network of [2] underperforms both EPN [3] and our network. One of the reasons could be that both EPN and our model for the classification task are also trained with the auxiliary task of rotation estimation (for the GA pooling layer or the permutation layer). In contrast, [2] is only trained with the classification task as described in their paper, which may cause more information loss in the invariant layer. Another reason could be that voxel inputs lose some details compared with point clouds. Transitioning from ModelNet10 to ModelNet40 dataset may require some scale-up of the network in terms of depth and width and some other careful tuning.

References

- [1] Tullio Ceccherini-Silberstein, A Machì, Fabio Scarabotti, and Filippo Tolli. Induced representations and Mackey theory. *Journal of Mathematical Sciences*, 156(1):11–28, 2009.
- [2] Gabriele Cesa, Leon Lang, and Maurice Weiler. A program to build $E(N)$ -equivariant steerable CNNs. In *International Conference on Learning Representations*, 2022.
- [3] Haiwei Chen, Shichen Liu, Weikai Chen, Hao Li, and Randall Hill. Equivariant point network for 3D point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 14514–14523, 2021.
- [4] Taco S Cohen, Mario Geiger, and Maurice Weiler. Intertwiners between induced representations (with applications to the theory of equivariant neural networks). *arXiv preprint arXiv:1803.10743*, 2018.
- [5] Taco S Cohen and Max Welling. Steerable cnns. *arXiv preprint arXiv:1612.08498*, 2016.
- [6] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [7] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotequi, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [8] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.