



Stony Brook University

# Frame Averaging

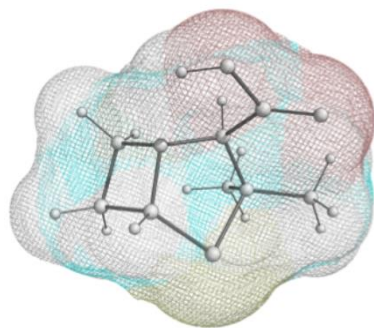
Wenhan Gao

Department of Applied Mathematics and Statistics

# Introduction

Group convolution is a popular framework for equivariant neural networks in fields such as computer vision, where the data can be treated as *functions*, and transformations act on the domain of these functions. In today's talk, we will delve into another framework, group averaging, for designing equivariant networks. Group averaging is currently mainly applied to "*geometric data*", e.g., point clouds or molecular structures that can be viewed as a set of "points/nodes" with coordinate information.

A Particular application is in geometric graph networks for computational physics, computational chemistry, 3D computer vision, etc.. It is desirable to have  $E(3)$  equivariant graph NNs in these tasks. Maintaining coordinate information offers *greater expressivity* compared to invariant methods that rely on relative distances [1].



The [pyridine](#) molecule can be described in the XYZ format by the following:

```
11
C      -0.180226841    0.360945118   -1.120304970
C      -0.180226841    1.559292118   -0.407860970
C      -0.180226841    1.503191118    0.986935030
N      -0.180226841    0.360945118    1.29018350
C      -0.180226841   -0.781300882    0.986935030
C      -0.180226841   -0.837401882   -0.407860970
H      -0.180226841    0.360945118   -2.206546970
H      -0.180226841    2.517950118   -0.917077970
H      -0.180226841    2.421289118    1.572099030
H      -0.180226841   -1.699398882    1.572099030
H      -0.180226841   -1.796059882   -0.917077970
```



## Definition: Group Averaging

Consider an arbitrary  $\Phi : X \rightarrow Y$ .

- $X, Y$ : Input and output (vector) spaces

The GA operator  $\langle \Phi \rangle_G : X \rightarrow Y$  is defined as:

$$\langle \Phi \rangle_G(x) = \mathbb{E}_{g \sim \nu} \rho_2(g) \cdot \Phi(\rho_1(g)^{-1} \cdot x) = \int_G \rho_2(g) \cdot \Phi(\rho_1(g)^{-1} \cdot x) d\nu(x)$$

or in summation form for discrete groups:

$$\langle \Phi \rangle_G(x) = \frac{1}{|G|} \sum_{g \in G} \rho_2(g) \cdot \Phi(\rho_1(g)^{-1} \cdot x)$$

- $\rho_1(g), \rho_2(g)$ : Group representations on  $X$  and  $Y$  respectively.
- $\nu$ : Harr measure over  $G$  ("uniform" over  $G$ )

Issue: When  $|G|$  is large (e.g., combinatorial groups such as permutations) or infinite (e.g., continuous groups such as rotations), then exact averaging is intractable.



# Intuition: Group Averaging

The GA operator is equivariant to  $G$ . Proof:

$$\begin{aligned}\langle \Phi \rangle_G(h \cdot x) &= \mathbb{E}_{g \sim \nu} \rho_2(g) \cdot \Phi(\rho_1(g)^{-1} \cdot (\rho_1(h) \cdot x)) \\ &= \mathbb{E}_{g \sim \nu} \rho_2(g) \cdot \Phi(\rho_1(h^{-1}g)^{-1} \cdot x) \\ &= \rho_2(h) \mathbb{E}_{g \sim \nu} \rho_2(h^{-1}g) \cdot \Phi(\rho_1(h^{-1}g)^{-1} \cdot x) \\ &= \rho_2(h) \langle \Phi \rangle_G(x)\end{aligned}$$

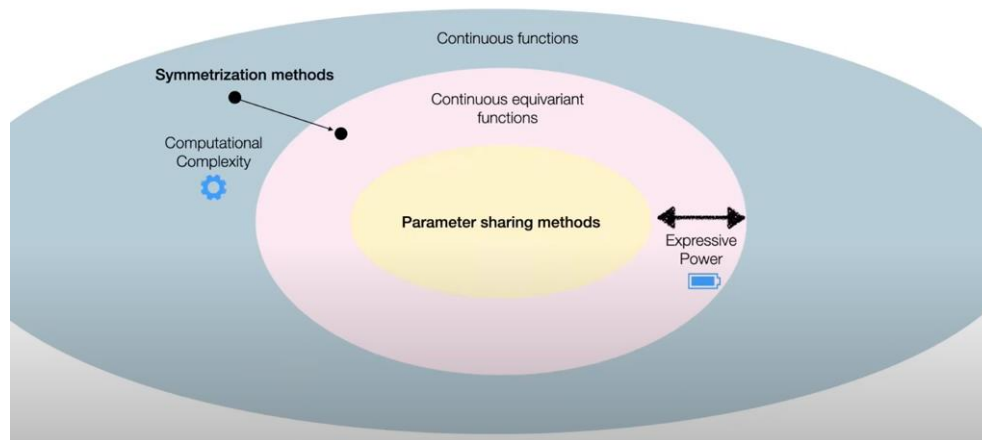
Intuition:

- Similar to group convolutions, we have already calculated all the transformed versions of the input  $(\rho_1(g)^{-1})$ .
- $\rho_2(g)$  "corrects" the output for equivariance.
  - $\{\Phi(\rho_1(g)^{-1} \cdot x), \forall g\}$  will result in the same set of outputs, but in a different order, for transformed inputs.
  - Why? Because the set of inputs  $\{\rho_1(g)^{-1} \cdot x\}$  is the same but in a different order for a transformed  $x$ .
  - Thus, integrating/summing over these outputs will result in invariant outputs.
  - $\rho_2(g)$  "corrects" the output by applying the transformation back.

# Expressive Power: Group Averaging

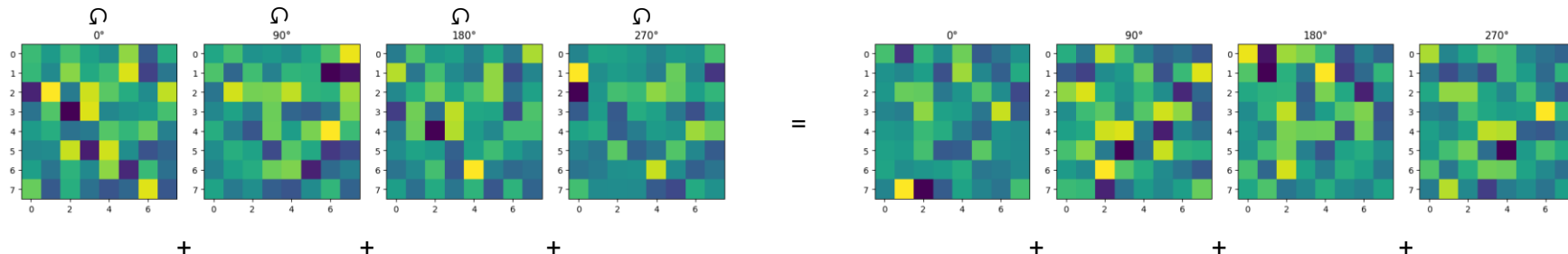
GA is as expressive as its backbone  $\Phi$  when  $\Phi$  is equivariant to  $G$ . To see this:

$$\begin{aligned}\langle \Phi \rangle_G(x) &= \mathbb{E}_{g \sim \nu} g \cdot \Phi(g^{-1} \cdot x) \\ &= \mathbb{E}_{g \sim \nu} g \cdot g^{-1} \cdot \Phi(x) \\ &= \Phi(x)\end{aligned}$$



# Connections between GA and G-Conv

Recall that in Group Conv, a rotation to the input will result in planar rotation + periodic shifting for the output feature map. If **we average over the feature maps to get a single feature map, the result is still equivariant to rotations, i.e., a rotation in the input will result in a rotation in the output feature map.**



Recall the definition of lifting layer:  $[f \star k](g) = \int_{\mathbb{R}^d} f(y)k(g^{-1}y) dy = (L_g \cdot k) \star f$ .

Make cross-correlation to be the backbone for group averaging:

$$\langle \Phi \rangle_G(f) = \mathbb{E}_{g \sim \nu} L_g \cdot [k \star (L_g^{-1} \cdot f)] = \frac{1}{|G|} \sum_{g \in G} L_g \cdot [k \star (L_g^{-1} \cdot f)] = \frac{1}{|G|} \sum_{g \in G} \overbrace{(L_g \cdot k) \star f}^{\text{Sum over all feature maps}}$$

Here, we have the “correction term”, but in group convolutions, we do not have the “correction term”. This is the result of moving the rotation from the image to the kernel. We’ll demonstrate this in the next slide.

A feature map after lifting



# Connections between GA and G-Conv

- Rotation to kernel instead of image:  $k \star L_g f = L_g (L_g^{-1} k \star f)$
- Intuition: Just distribute  $L_g$  inside (However, intuition is wrong here, mathematically this is incorrect, and for many representations, this is wrong, we can do this here because rotations are unitary matrices).
- Proof:

$$\begin{aligned}
 & [(L_g k) \star (L_g f)](x) \\
 &= \int_{\mathbb{R}^2} L_g k(x - y) L_g f(y) dy \\
 &= \int_{\mathbb{R}^2} k(gx - gy) f(gy) dy \\
 &= \int_{\mathbb{R}^2} k(gx - y') f(y') dy' \\
 &= L_g(k \star f)(x)
 \end{aligned}$$

- First equality: Just definition of conv/cross-correlation
- Second equality: Again just distribute  $L_g$
- Third equality: Change of variable  $gy = y'$  since the determinant of  $g$  is 1.
- Fourth equality: Definition



# Motivating Example: Pre-processing

Consider image segmentation, in which we want translation equivariance. Assuming we are using a model other than CNN, so that we do not have inherent translation equivariance. What can we do if we want translation equivariance? We can use group averaging, but it will be computationally intractable for the translation group (large in the discrete case or even infinite if we view it in a continuous manner). We should think of another way of achieving equivariance.

One way is geometric pre-processing:

Given an image,



we can achieve equivariance by preprocessing the image.

For example, if we have the location of the left eye of a cat, we can preprocess the image such that all cats will have their left eyes in the same location.



# Definition: Frame Averaging

A frame is defined as a set valued function  $\mathcal{F} : X \rightarrow 2^G$ . (Taking an input  $x \in X$  and mapping to a subset of  $G$ )

A frame is equivariant to  $G$  if  $\mathcal{F}(g \cdot x) = g\mathcal{F}(x), \forall g \in G$ . (equality of sets)

Frame Averaging operator is defined as:

$$\langle \Phi \rangle_{\mathcal{F}}(x) = \frac{1}{|\mathcal{F}(x)|} \sum_{g \in \mathcal{F}(x)} \rho_2(g) \Phi(\rho_1(g)^{-1}x)$$

Similar to group averaging, we can prove that frame averaging operator is equivariant to  $G$  if  $\mathcal{F}$  is equivariant to  $G$ , and it is as expressive as its backbone if its backbone is equivariant.

Example:

Consider  $X = \mathbb{R}^n, Y = \mathbb{R}^n$ , and  $G = \mathbb{R}$  with addition as the group action. We choose the group actions in this case to be  $\rho_1(t)\mathbf{x} = \mathbf{x} + t\mathbf{1}$ , and  $\rho_2(a)y = y + t$ , where  $t \in G, \mathbf{x} \in X, y \in Y$ , and  $\mathbf{1} \in \mathbb{R}^n$  is the vector of all ones.

We can define the frame in this case using the averaging operator  $\mathcal{F}(\mathbf{x}) = \{\frac{1}{n}\mathbf{1}^T \mathbf{x}\} \subset G = \mathbb{R}$ .

Note that in this case the frame contains only one element from the group, in other cases finding such a small frame is hard or even impossible.

One can check that this frame is equivariant. The FA:  $\langle \Phi \rangle_{\mathcal{F}}(\mathbf{x}) = \Phi\left(\mathbf{x} - \frac{1}{n}(\mathbf{1}^T \mathbf{x})\mathbf{1}\right) + \frac{1}{n}\mathbf{1}^T \mathbf{x}$  in the equivariant case.

- Intuition: Geometric pre-processing, we subtract the average and then add the average back to obtain equivariance.

# Frame Averaging: Practical Usage in Geometric GNNs

Let  $\Phi : V \rightarrow W$  be an arbitrary function, where  $V, W$  are some vector spaces. The group averaging operator  $\Psi$  can be made equivariant by symmetrization, that is averaging over the group:

$$\Psi(X) = \frac{1}{|G|} \sum_{g \in G} g \cdot \Phi(g^{-1} \cdot X).$$

- $\Phi$ : Node update MLPs in GNN
- $X \in V$ : Input features (coordinates in particular)
- $W$ : Space of output node embeddings
- $G, |G|$ : Group and cardinality of the group, respectively

It can be shown that  $\Psi : V \rightarrow W$  is equivariant w.r.t.  $G$ .

Or alternatively, if we can find a frame for the group, we can use frame averaging instead.



# Frame Averaging: Practical Instantiation for E(3)

Goal:

- We would like to incorporate Euclidean symmetry to existing permutation invariant/equivariant point cloud networks.

Settings:

- Input Space:  $V = \mathbb{R}^{n \times d}$  ( $n$  nodes, each holding a  $d$ -dimensional vector as its location)
- Group:  $G = E(d) = O(d) \times T(d)$ , namely the group of Euclidean motions in  $\mathbb{R}^d$  defined by rotations and reflections  $O(d)$ , and translations  $T(d)$ .
- Representation acting on  $\mathbf{X} \in V$ :  $\rho_1(g)\mathbf{X} = \mathbf{X}\mathbf{R}^T + \mathbf{1}\mathbf{t}^T$ , where  $\mathbf{R} \in O(d)$ , and  $\mathbf{t} \in \mathbb{R}^d$ . (Apply rotation and translation to every node)
- $W, \rho_2$  are defined similarly, unless we want invariance.

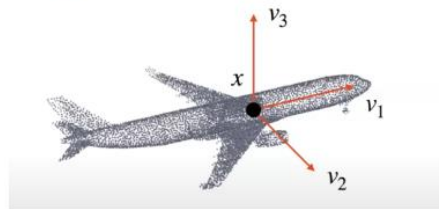


# Frame Averaging: Practical Instantiation for E(3)

Frame  $\mathcal{F}(\mathbf{X})$  is defined based on Principle Component Analysis (PCA), as follows:

- Let  $\mathbf{t} = \frac{1}{n} \mathbf{X}^T \mathbf{1} \in \mathbb{R}^d$  be the centroid of  $\mathbf{X}$
- $\mathbf{C} = (\mathbf{X} - \mathbf{1t}^T)^T (\mathbf{X} - \mathbf{1t}^T) \in \mathbb{R}^{d \times d}$  the covariance matrix computed after removing the centroid from  $\mathbf{X}$ . In the generic case the eigenvalues of  $\mathbf{C}$  satisfy  $\lambda_1 < \lambda_2 < \dots < \lambda_d$ .
- Let  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$  be the unit length corresponding eigenvectors.
- Then we define  $\mathcal{F}(\mathbf{X}) = \{([\alpha_1 \mathbf{v}_1, \dots, \alpha_d \mathbf{v}_d], \mathbf{t}) \mid \alpha_i \in \{-1, 1\}\} \subset E(d)$ .
- $[\mathbf{v}_1, \dots, \mathbf{v}_d]$  is a set of orthonormal vectors in  $\mathbb{R}^d$ , i.e., a basis of  $\mathbb{R}^d$ . Moreover, these vectors will "rotate" in the same way as the input.
- $\mathcal{F}(\mathbf{X})$  based on the covariance and centroid are  $E(d)$  equivariant.

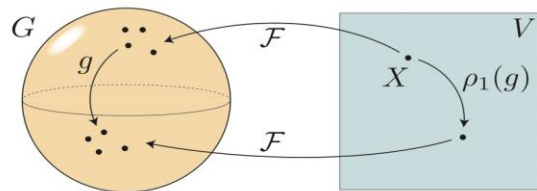
Example in 3D:



# Intuition: Frame Averaging

Overall, frame averaging achieves equivariance by:

1. Create a "coordinate system" (or a set of orthonormal basis in the general case) that will change according to the group actions acting on the input.
  - Meaning that if one input is a transformed version of another, this coordinate system will be transformed in a predictable way!
2. Representing the input in terms of this "coordinate system".
  - Inputs will be the same after this step (if one is a transformed version of another), the group action is reflected on the coordinate system now!
3. Reconstruction by the "coordinate system".
  - Now the group action is reflected on the output!



$$\langle \Phi \rangle_{\mathcal{F}}(x) = \frac{1}{|\mathcal{F}(x)|} \sum_{g \in \mathcal{F}(x)} \rho_2(g) \Phi(\rho_1(g)^{-1} x)$$



Mapping the input to a set of "coordinate systems" that respects transformations on the input.



Just synchronize over all possible such "coordinate systems".



# Questions for Averaging-based Methods

- ❑ Most applications of frame averaging are applied to geometric data (point clouds). Can we adapt this idea for other data types, e.g., image data, functional data?
- ❑ Can we have an averaging-based design that also respects local symmetries?
- ❑ Under what conditions can we find a frame, and under what conditions can we find an optimal frame?
- ❑ ...

