

**Sinkformers: Transformers with Doubly Stochastic Attention  
(AISTATS2022)**

Michael E. Sander, Pierre Ablin.etc.

**LieTransformer: Equivariant Self-Attention for Lie Groups  
(ICML2021)**

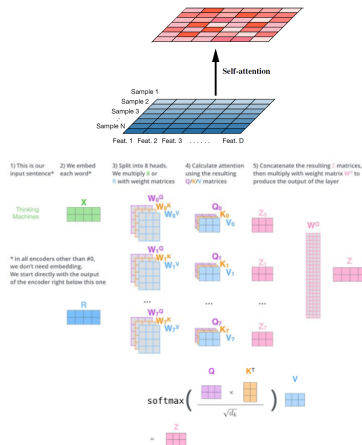
Michael Hutchinson, Charline Le Lan.etc.

**Presenter: Minjie Cheng**

# Outline

1. Self-Attention
2. Sinkformers
  - 2.1. Motivation
  - 2.2. Algorithm
  - 2.3. Experiments
3. LieTransformer
  - 3.1. Motivation
  - 3.2. Algorithm
  - 3.3. Experiments

# Overview of Self-Attention



## ► A mapping

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{N \times D} \mapsto \mathbf{Z} \in \mathbb{R}^{N \times D}$$

## ► Using $m$ to index the head

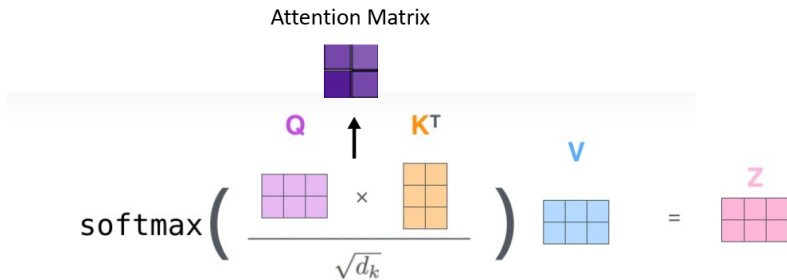
( $m = 1, \dots, M$ ), the output of the  $m$ th head can be written as:

$$f^m(\mathbf{X}) \triangleq \mathbf{X} \mathbf{W}^{\mathbf{Q},m} (\mathbf{X} \mathbf{W}^{\mathbf{K},m})^T \in \mathbb{R}^{N \times D/M}$$

$$\mathbf{W} \triangleq \text{softmax}(\mathbf{X} \mathbf{W}^{\mathbf{Q},m} (\mathbf{X} \mathbf{W}^{\mathbf{K},m})^T) \in \mathbb{R}^{N \times N}$$

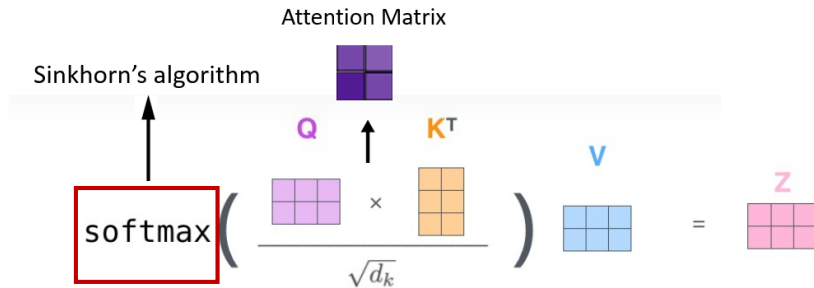
$$\text{MSA}(\mathbf{X}) \triangleq [f^1(\mathbf{X}), \dots, f^M(\mathbf{X})] \mathbf{W}^{\mathbf{O}} \in \mathbb{R}^{N \times D}$$

# Motivation-Sinkformers



- This attention matrix is normalized with the SoftMax operator, which makes it row-wise stochastic

# Motivation-Sinkformers

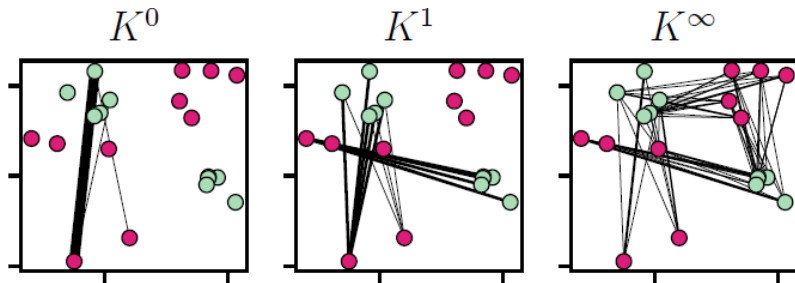


- ▶ This attention matrix is normalized with the SoftMax operator, which makes it row-wise stochastic
- ▶ Use Sinkhorn's algorithm to make attention matrices **doubly** stochastic.(i.e., rows and columns both sum to 1)

# Important Notations

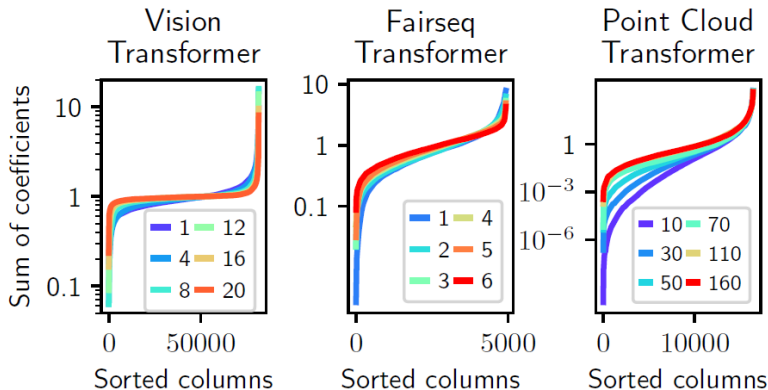
- ▶  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{N \times D}$
- ▶  $\mathbf{x}_i \leftarrow \mathbf{x}_i + \sum_{j=1}^n \mathbf{K}_{ij}^1 \mathbf{W}_V \mathbf{x}_j$ 
  - ▶  $\mathbf{K}^1 := \text{Softmax}(\mathbf{C})$
  - ▶  $C_{i,j} = (\mathbf{W}_Q \mathbf{x}_i)^T \mathbf{W}_K \mathbf{x}_j$
  - ▶ Query, key, value matrices are  $\mathbf{W}_Q \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_K \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_V \in \mathbb{R}^{d \times d}$
  - ▶ Denote  $\mathbf{K}^0 := \exp(\mathbf{C})$
  - ▶  $K_{ij}^1 := K_{ij}^0 / \sum_{l=1}^n K_{il}^0$
  - ▶  $\mathbf{K}^1$  is row-wise stochastic

# Illustration of the different normalizations of attention matrices



- ▶  $C_{i,j} = (\mathbf{W}_Q \mathbf{x}_i)^T \mathbf{W}_K \mathbf{x}_j$  (correlation coefficient)
- ▶ We form two point clouds  $\mathbf{W}_Q \mathbf{x}_i$  (green) and  $\mathbf{W}_K \mathbf{x}_j$  (red).
- ▶ We only display connections with  $K_{ij}^k \geq 10^{-12}$
- ▶ For  $K^\infty$  (Sinkhorn), all points are involved in an interaction.

## Experiment for Sum over columns



- ▶ We sum over columns of attention matrices at different training epochs (color) when training.
- ▶ The majority of columns naturally sum closely to 1.



# Sinkhorn's algorithm

$$K^{l+1} = \begin{cases} N_R(K^l) & \text{if } l \text{ is even} \\ N_C(K^l) & \text{if } l \text{ is odd} \end{cases}$$

- $N_R$  and  $N_C$  correspond to row-wise and column-wise normalizations:

$$(N_R(K))_{ij} := \frac{K_{ij}}{\sum_{l=1}^n K_{i,l}}$$

$$(N_C(K))_{ij} := \frac{K_{ij}}{\sum_{l=1}^n K_{l,j}}$$

- Note that it is doubly stochastic in the sense that  $K^\infty \mathbb{1}_n = \mathbb{1}_n$  and  $K^{\infty T} \mathbb{1}_n = \mathbb{1}_n$ .

# Sinkformers

- Note that  $K^1 := \text{SoftMax}(C)$  is precisely the output of Sinkhorn's algorithm after 1 iteration.

$$x_i \leftarrow x_i + \sum_{j=1}^n K_{i,j}^{\infty} W_V x_j$$

- Only a few iterations of Sinkhorn are sufficient (**typically 3 to 5**) to converge to a doubly stochastic matrix.
- The practical training time of Sinkformers is comparable to regular Transformers.

# Invariance to the cost function

**Proposition 1.** Let  $C \in \mathbb{R}^{n \times n}$ . Consider, for  $(f, g) \in \mathbb{R}^n \times \mathbb{R}^n$  the modified cost function  $\tilde{C}_{ij} := C_{ij} + f_i + g_j$ . Then  $\text{Sinkhorn}(C) = \text{Sinkhorn}(\tilde{C})$ .

$$\tilde{C}_{ij} := -\frac{1}{2} \|W_Q x_i - W_K x_j\|^2$$

- We can consider the cost  $\tilde{C}_{ij}$  instead of  $C_{ij} = (W_Q x_i)^\top W_K x_j$ , without affecting  $K^\infty$ .

## Residual maps for attention-continuous counterparts

- We denote  $c(x, x') := (W_Q x)^\top W_K x'$  and  $k^0 := \exp(c)$ . For some measure  $\mu \in \mathcal{M}(\mathbb{R}^d)$ , we define the SoftMax operator on the cost  $c$  by  $k^1(x, x') = \text{SoftMax}(c)(x, x') := \frac{k^0(x, x')}{\int k^0(x, y) d\mu(y)}$ . Similarly, we define Sinkhorn's algorithm as the following iterations, starting from  $k^0 = \exp(c)$  :

$$k^{l+1}(x, x') = \begin{cases} \frac{k^l(x, x')}{\int k^l(x, y) d\mu(y)} & \text{if } l \text{ is even} \\ \frac{k^l(x, x')}{\int k^l(y, x) d\mu(y)} & \text{if } l \text{ is odd.} \end{cases}$$

We denote  $k^\infty := \text{Sinkhorn}(c)$  the resulting limit. Note that if  $\mu$  is a discrete measure supported on a  $n$  sequence of particles  $(x_1, x_2, \dots, x_n)$ ,  $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ , then for all  $(i, j)$ ,  $k^0(x_i, x_j) = K_{i,j}^0$ ,  $k^1(x_i, x_j) = K_{i,j}^1$  and  $k^\infty(x_i, x_j) = K_{i,j}^\infty$ , so that  $k^0, k^1$  and  $k^\infty$  are indeed the **continuous equivalent** of the matrices  $K^0, K^1$  and  $K^\infty$  respectively.

# depth limit

- ▶ Transformer equation

$$x_i \leftarrow x_i + \sum_{j=1}^n K_{ij}^1 W_V x_j$$

- ▶ ResNet equation

$$x_i \leftarrow x_i + T(x_i)$$

- ▶ Sinkformer equation

$$x_i \leftarrow x_i + \sum_{j=1}^n K_{ij}^\infty W_V x_j$$

## Infinitesimal step-size regime

- ▶ In this framework, iterating the Transformer equation, the ResNet equation and the Sinkformer equation corresponds to a Euler discretization with step-size 1 of the ODEs

$$\dot{x}_i = T_\mu(x_i) \text{ for all } i,$$

- ▶ where  $x_i(t)$  is the position of  $x_i$  at time  $t$ . For an arbitrary measure  $\mu \in \mathcal{M}(\mathbb{R}^d)$ , these ODEs can be equivalently written as a continuity equation

$$\partial_t \mu + \operatorname{div}(\mu T_\mu) = 0$$

- ▶ When  $T_\mu$  is defined by the ResNet equation 2,  $T_\mu = T$  does not depend on  $\mu$ . It defines an advection equation where the particles do not interact and evolve independently.
- ▶  $T_\mu^1(x) = \int k^1(x, x') W_{V'} x' d\mu(x')$
- ▶  $T_\mu^\infty(x) = \int k^\infty(x, x') W_{V'} x' d\mu(x')$

## Implementation details

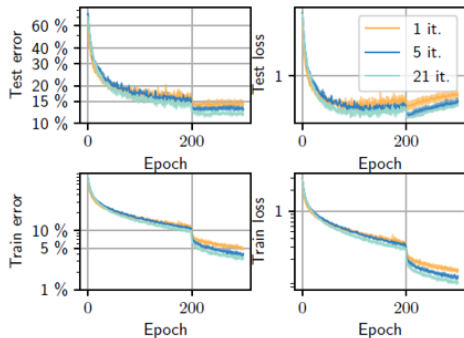
- We implement Sinkhorn's algorithm in log domain for stability. Given a matrix  $K^0 \in \mathbb{R}^{n \times n}$  such that  $K_{ij}^0 = e^{C_{ij}}$  for some  $C \in \mathbb{R}^{n \times n}$ , Sinkhorn's algorithm approaches  $(f, g) \in \mathbb{R}^n \times \mathbb{R}^n$  such that  $K^\infty = \text{diag}(e^{f^\infty}) K^0 \text{diag}(e^{g^\infty})$  by iterating in log domain, starting from  $g^0 = \mathbb{O}_n$ ,

$$\begin{aligned} f^{l+1} &= \log(\mathbb{1}_n/n) - \log(Ke^{g^l}) && \text{if } l \text{ is even} \\ g^{l+1} &= \log(\mathbb{1}_n/n) - \log(K^T e^{f^l}) && \text{if } l \text{ is odd.} \end{aligned}$$

This allows for fast and accurate computations, where  $\log(Ke^{g^l})$  and  $\log(K^T e^{f^l})$  are computed using log-sum-exp.

# Experiments

## ► ModelNet 40 classification



Model	Best	Median	Mean	Worst
Set Transformer	87.8%	86.3%	85.8%	84.7%
Set Sinkformer	89.1%	88.4%	88.3%	88.1%
Point Cloud Transformer	93.2%	92.5%	92.5%	92.3%
Point Cloud Sinkformer	93.1%	92.8%	92.7%	92.5%



# Experiments

## ► Sentiment Analysis

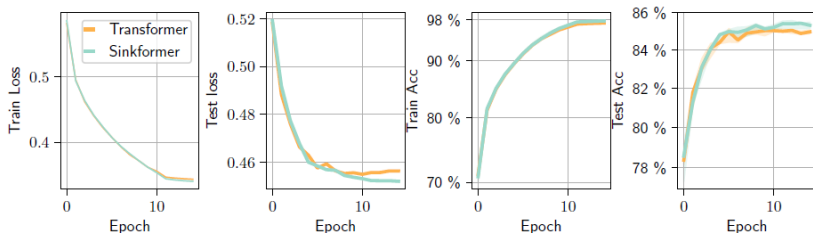
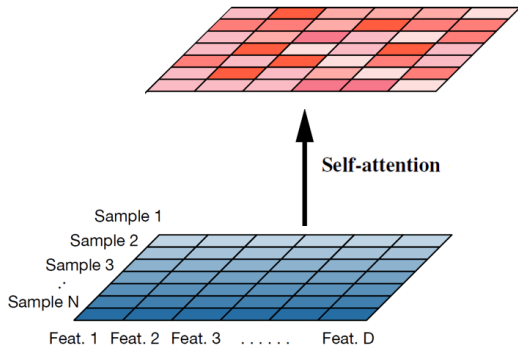


Figure 4: Learning curves when training a Transformer and a Sinkformer on the Sentiment Analysis task on the IMDb Dataset.

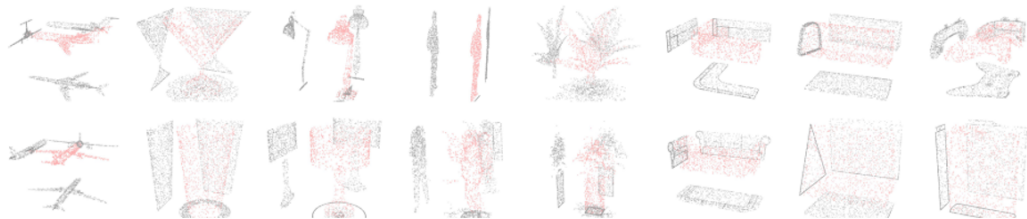
## **2. LieTransformer: Equivariant Self-Attention for Lie Groups** (ICML2021)

# Motivation-LieTransformer



- Extend **group equivariance** to self-attention, **non-linear** map

# Motivation-LieTransformer



- ▶ 3D point clouds
- ▶ We may want a classifier to output the same classification when the input is translated or rotated.
- ▶ Transformer-based model that is **permutation invariant**, but **not invariant to rotations or translations**.

# Equivariant Self-Attention for Lie Groups

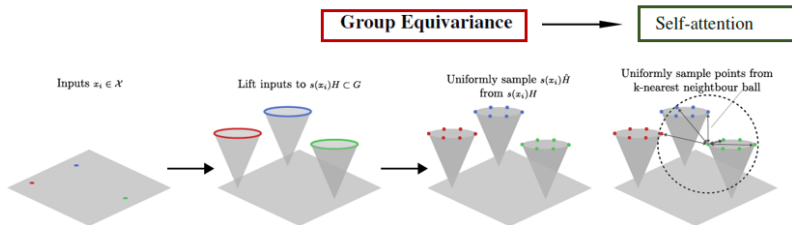


Figure 2. Visualisation of lifting, sampling  $\hat{H}$ , and subsampling in the local neighbourhood for  $SE(2)$  acting on  $\mathbb{R}^2$ . Self-attention is performed on this subsampled neighbourhood.

# Formal definitions for Groups and Representation Theory

- ▶ **Group**  $G$  is a set of symmetries, with each group element  $g$  corresponding to a symmetry transformation, group element  $g \in G$
- ▶ **Graph representation**  $\rho(g)$  takes the form of a matrix. For  $SO(2)$ ,

$$\rho(g_\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

- ▶ **Definition.** A group  $G$  is a set endowed with a single operator :  $G \times G \mapsto G$  such that
  1. Associativity:  $\forall g, g', g'' \in G, (g \cdot g') \cdot g'' = g \cdot (g' \cdot g'')$
  2. Identity:  $\exists e \in G, \forall g \in G, g \cdot e = e \cdot g = g$
  3. Invertibility:  $\forall g \in G, \exists g^{-1} \in G, g \cdot g^{-1} = g^{-1} \cdot g = e$

# $G$ -equivariance

- ▶ Definition 1. We say that a map  $\Phi : V_1 \rightarrow V_2$  is  $G$  equivariant with respect to actions  $\rho_1, \rho_2$  of  $G$  acting on  $V_1, V_2$  respectively if:  $\Phi [\rho_1(g)f] = \rho_2(g)\Phi[f]$  for any  $g \in G, f \in V_1$ .
- ▶ In the context of group equivariant neural networks,  $V$  is commonly defined to be the space of scalar-valued functions on some set  $S$ , so that  $V = \{f \mid f : S \rightarrow \mathbb{R}\}$ .
- ▶ i.e. a grey-scale image can be expressed as a feature map  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  from pixel coordinate  $x_i$  to pixel intensity  $f_i$ , supported on the grid of pixel coordinates. (Linear operation )

# Overview

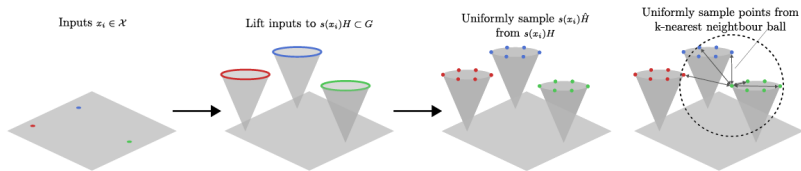


Figure 2. Visualisation of lifting, sampling  $\hat{H}$ , and subsampling in the local neighbourhood for  $SE(2)$  acting on  $\mathbb{R}^2$ . Self-attention is performed on this subsampled neighbourhood.



# Lifting

- ▶ Suppose we have data in the form of a set of input pairs  $(x_i, \mathbf{f}_i)_{i=1}^n$  where  $x_i \in \mathcal{X}$  are **spatial coordinates** and  $\mathbf{f}_i \in \mathcal{F}$  are **feature values**.
- ▶ All elements of  $\mathcal{X}$  are connected by the action:  $\forall x, x' \in \mathcal{X}, \exists g \in G : \rho(g)x = x'$

# LieSelfAttention

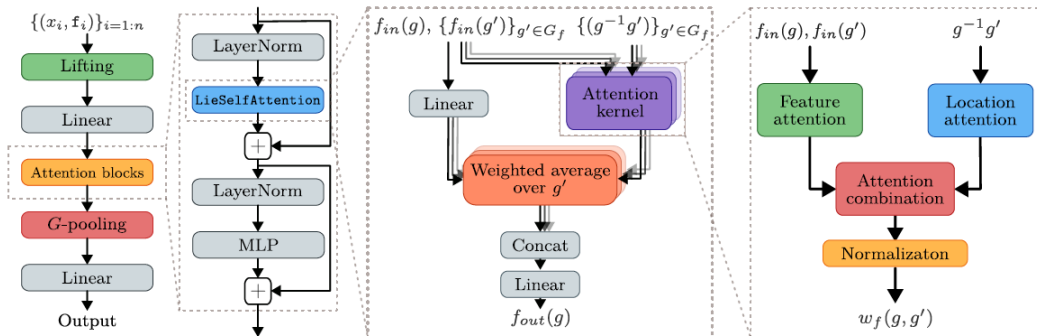


Figure 1. Architecture of the LieTransformer.

# Content-based attention

- Content-based attention  $k_c(f(g), f(g'))$  :
  1. Dot-product:  $\frac{1}{\sqrt{d_v}} (W^Q f(g))^\top W^K f(g') \in \mathbb{R}$  for  $W^Q, W^K \in \mathbb{R}^{d_v \times d_v}$
  2. Concat:  $\text{Concat} [W^Q f(g), W^K f(g')] \in \mathbb{R}^{2d_v}$
  3. Linear-Concat-linear:  $W \text{Concat} [W^Q f(g), W^K f(g')] \in \mathbb{R}^{d_s}$  for  $W \in \mathbb{R}^{d_s \times 2d_v}$ .

# Location-based attention

- ▶ Location-based attention  $k_l(g^{-1}g')$  for Lie groups  $G$  :
  1. Plain:  $\nu [\log (g^{-1}g')]$
  2. MLP:  $\text{MLP} (\nu [\log (g^{-1}g')])$

# Experiments-Counting Shapes in 2D Point Clouds

Training data	$D_{\text{train}}$	$D_{\text{train}}$	$D_{\text{train}}$	$D_{\text{train}}^{T^2}$	$D_{\text{train}}^{T^2}$	$D_{\text{train}}^{SE^2}$
Test data	$D_{\text{test}}$	$D_{\text{test}}^{T^2}$	$D_{\text{test}}^{SE^2}$	$D_{\text{test}}^{T^2}$	$D_{\text{test}}^{SE^2}$	$D_{\text{test}}^{SE^2}$
SetTransformer	$0.58 \pm 0.07$	$0.44 \pm 0.02$	$0.44 \pm 0.02$	$0.61 \pm 0.02$	$0.51 \pm 0.01$	$0.55 \pm 0.01$
LieTransformer-T2	$0.75 \pm 0.03$	$0.75 \pm 0.03$	$0.63 \pm 0.06$	$0.75 \pm 0.03$	$0.63 \pm 0.06$	$0.70 \pm 0.03$
LieTransformer-SE2	$0.71 \pm 0.01$	$0.71 \pm 0.01$	$0.69 \pm 0.02$	$0.71 \pm 0.01$	$0.69 \pm 0.02$	$0.72 \pm 0.04$

Table 1. Mean and standard deviation of test accuracies on the shape counting task at convergence (over 8 random initialisations).

- ▶ Transformer-based model that is **permutation invariant**, but **not invariant to rotations or translations**.
- ▶  $G$ -equivariance: rotations, translation

# Experiments-QM9: Molecule Property Regression

Task Units	$\alpha$ bohr <sup>3</sup>	$\Delta\epsilon$ meV	$\epsilon_{\text{HOMO}}$ meV	$\epsilon_{\text{LUMO}}$ meV	$\mu$ D	$C_v$ cal/mol K	$G$ meV	$H$ meV	$R^2$ bohr <sup>2</sup>	$U$ meV	$U_0$ meV	ZPVE meV
WaveScatt (Hirn et al., 2017)	.160	118	85	76	.340	.049	—	—	—	—	—	—
NMP (Gilmer et al., 2017)	<b>.092</b>	69	43	38	<b>.030</b>	.040	19	17	.180	20	20	<b>1.50</b>
SchNet (Schütt et al., 2017)	.235	<b>63</b>	<b>41</b>	<b>34</b>	.033	<b>.033</b>	<b>14</b>	<b>14</b>	<b>.073</b>	<b>19</b>	<b>14</b>	1.70
Cormorant (Anderson et al., 2019)	.085	61	34	38	.038	.026	20	21	.961	21	22	2.03
DimeNet++ (Klicpera et al., 2020) *	.049	34	26	20	.033	.024	7.7	7.1	.387	6.7	6.9	1.23
LINet (Miller et al., 2020)	.088	68	45	35	.043	.031	14	14	.354	14	13	1.56
TFN (Thomas et al., 2018)	.223	58	40	38	.064	.101	—	—	—	—	—	—
SE3-Transformer (Fuchs et al., 2020)	.148	53	36	33	.053	.057	—	—	—	—	—	—
LieConv-T3 (Finzi et al., 2020) <sup>†</sup>	.125	60	36	32	.057	.046	35	37	1.54	36	35	3.62
LieConv-T3 + SO3 Aug (Finzi et al., 2020)	.084	49	30	<b>25</b>	<b>.032</b>	.038	22	24	.800	19	19	2.28
LieConv-SE3 (Finzi et al., 2020) <sup>†</sup>	.097	<b>45</b>	<b>27</b>	<b>25</b>	.039	.041	39	46	2.18	49	48	3.27
LieConv-SE3 + SO3 Aug (Finzi et al., 2020) <sup>†</sup>	.088	<b>45</b>	<b>27</b>	<b>25</b>	.038	.043	47	46	2.12	44	45	3.25
LieTransformer-T3 (Us)	.179	67	47	37	.063	.046	27	29	.717	27	28	2.75
LieTransformer-T3 + SO3 Aug (Us)	<b>.082</b>	51	33	27	.041	<b>.035</b>	<b>19</b>	<b>17</b>	<b>.448</b>	<b>16</b>	<b>17</b>	<b>2.10</b>
LieTransformer-SE3 (Us)	.104	52	33	29	.061	.041	23	27	2.29	26	26	3.55
LieTransformer-SE3 + SO3 Aug (Us)	.105	52	33	29	.062	.041	22	25	2.31	24	25	3.67

- **non-invariant** models specifically designed for the QM9 task.
- **invariant models** specifically designed for the QM9 task.

# Summary and Future Work

- ▶ Mapping
- ▶ UOT-based module replace the self-attention modules in Transformer.

