



香港中文大學

The Chinese University of Hong Kong



Ron

Sampling 3D Molecular Conformers with Diffusion Transformers

11 August 2025

Contents

- Introduction
- Results
- Technical Details
 - Background
 - Flow matching
 - $SO(3)$ -equivariant self-attention
 - Combination of methods in this work

Introduction

- Conformer generation with Diffusion model
- Less computationally expensive than DFT
- A small model in general (<100M parameters)

Results

Benchmark of this model against existing models

| Precision | Recall | MAE of ensemble properties | Out-of-distribution generalization |
|--|---|--|---|
| SOTA | On par | SOTA | SOTA |
| The least amount of incorrect samples generated. | The most amount of samples within the training set generated. | Energy, dipole moment, HOMO-LUMO gap, minimum energy | Testing with larger molecules (>100 atoms) with models trained on molecules with ~44 atoms on average |

Technical details

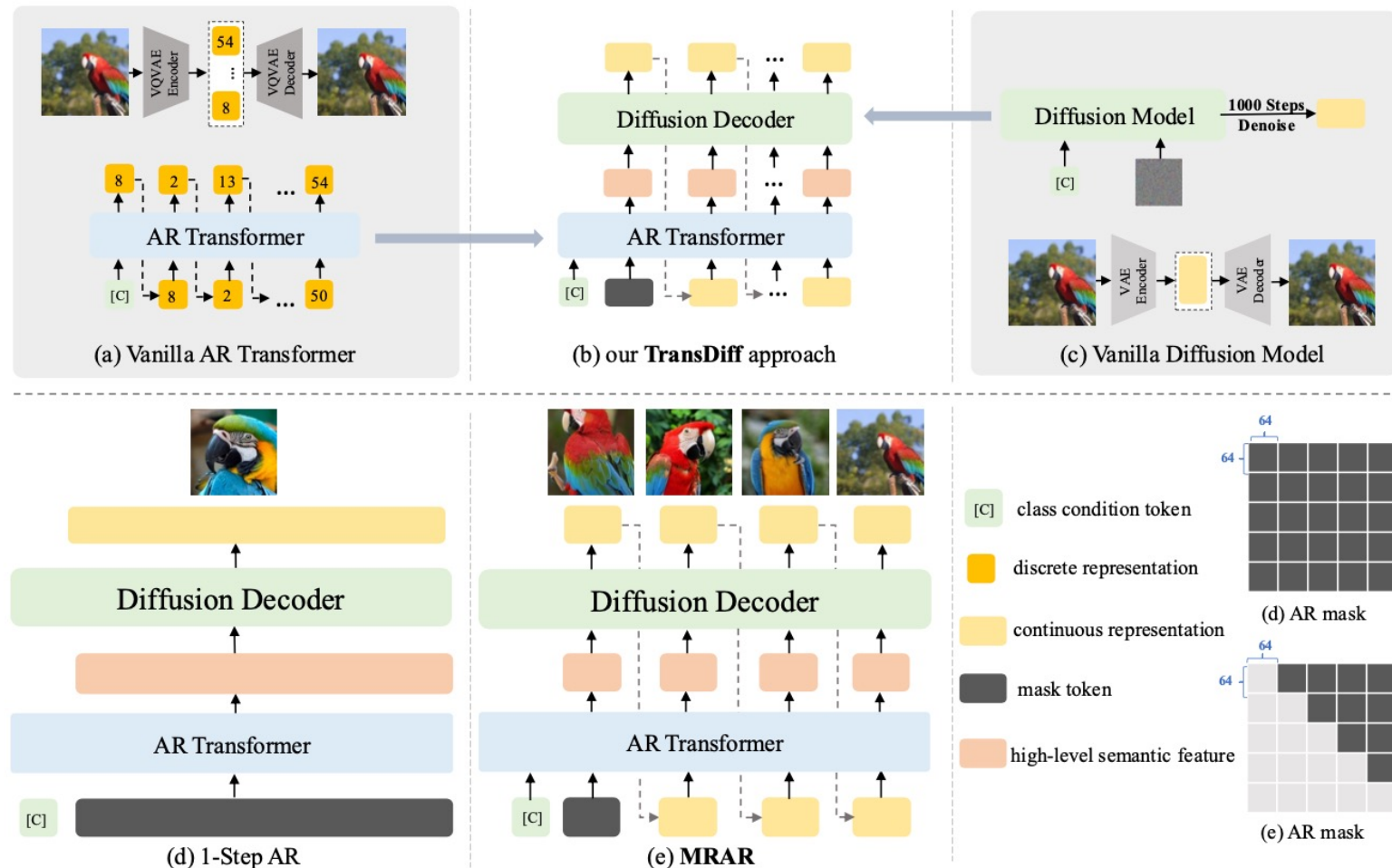
- Generative models
- Diffusion Transformers
- Flow matching
- Equivariance
- Conditioning
- Positional Embedding
- Self-attention

Generative models

A generative model is a joint probability distribution $p(x)$, for $x \in X$. In some cases, the model may be conditioned on inputs or covariates $c \in \mathcal{C}$, which gives rise to a conditional generative model of the form $p(x|C)$.

- Probabilistic Machine Learning: Advanced Topics

Diffusion Transformers



Equivariance

That the output of a function is identical for different inputs that can be transformed into each other via certain permitted operations.

Slides from
https://geometricdeeplearning.com/slides/Cambridge_1_Introduction_to_Groups_and_Representations.pdf

11 August 2025

Group invariance

We can now *formally* describe how to *exploit* the symmetries in \mathfrak{G} !

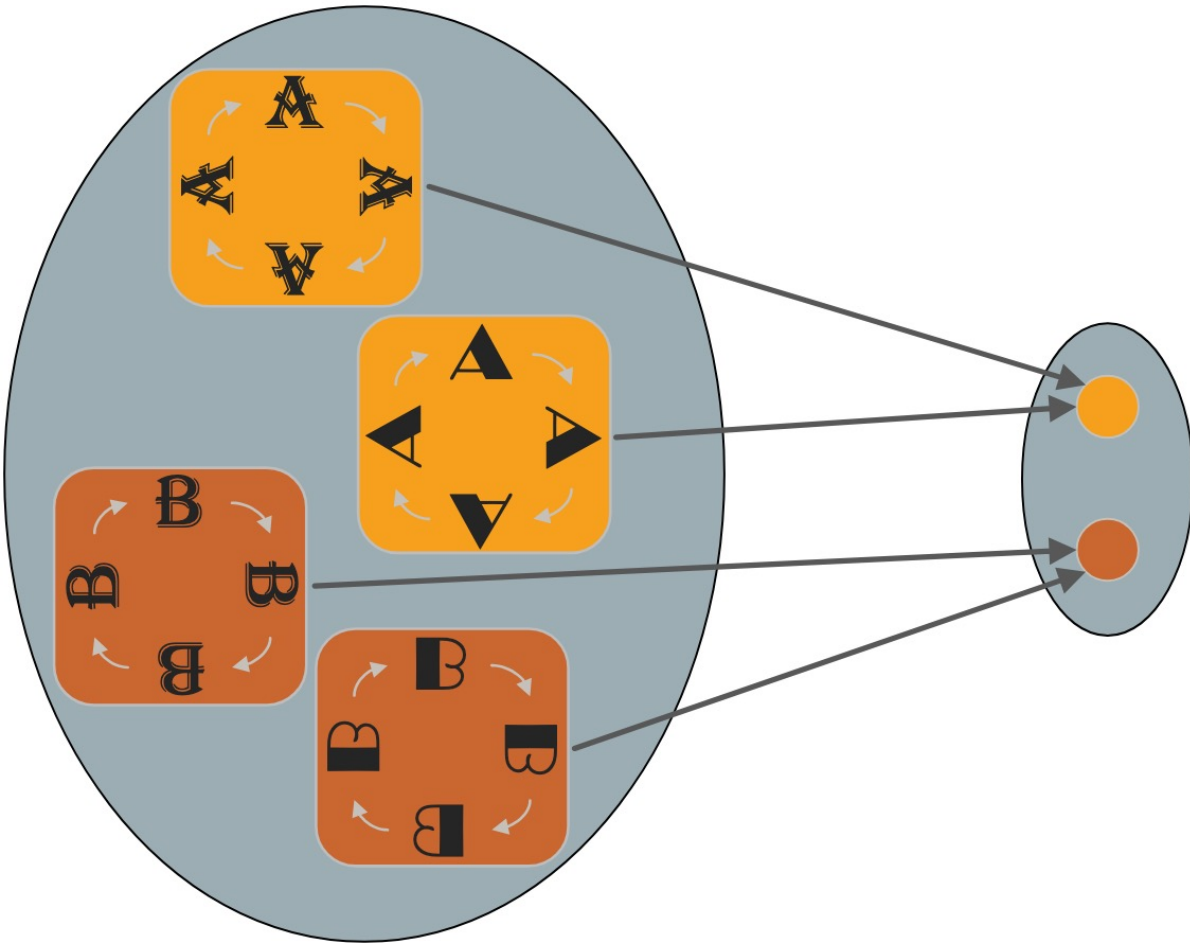
A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{Y}$ is \mathfrak{G} -invariant if $f(\rho_x(g)x) = f(x)$ for all $g \in \mathfrak{G}$, i.e., its output is unaffected by the group action on the input.

e.g. **image classification**: output class won't depend on image **shifts**

$$f\left(\begin{array}{c} \text{orange square} \\ \text{with beetle in bottom-right} \end{array}\right) = f\left(\begin{array}{c} \text{orange square} \\ \text{with beetle in top-left} \end{array}\right) = \text{beetle icon}$$

$\rho_x(g)x$ x

Orbits and equivalence relations



$$O_x = \{gx \mid x \in X, g \in G\}$$

G -equivalence

$$x \sim_G y \iff \exists g \in G : gx = y$$

Satisfies the axioms of an equivalence relation:

1. Reflexivity: $x \sim_G x$
 - (Because G contains the identity)
2. Transitivity:
$$x \sim_G y \wedge y \sim_G z \iff x \sim_G z$$
 - (Because G is closed under composition)
3. Symmetry: $x \sim_G y \iff y \sim_G x$
 - (Because G is closed under inverses)

Group equivariance

We proceed to define a more *fine-grained* notion of regularity:

A function $f : \mathcal{X}(\Omega) \rightarrow \mathcal{Z}(\Omega)$ is \mathfrak{G} -equivariant if, for all $g \in \mathfrak{G}$, $f(\rho_{\mathcal{X}}(g)x) = \rho_{\mathcal{Z}}(g)f(x)$, i.e., applying a group action on the input affects the output in the same way.

e.g. **image segmentation**: segmentation mask must **follow** any shifts in the input



Note that invariance is a *special case* of equivariance (for which $\rho_{\mathcal{Z}}$?)

Flow matching

We describe this transformation in terms of a *stochastic interpolant* \mathbf{x}_t [26, 49, 50]. A noisy sample at time $t \in [0, 1]$ is defined as

$$\mathbf{x}_t = (1 - t) \cdot \mathbf{x}_0 + t \cdot \mathbf{x}_1 + \sigma \cdot \epsilon, \quad (1)$$

where $\epsilon \in \mathbb{R}^{N \times 3}$ is drawn from the standard normal distribution $\mathcal{N}(0, \mathbf{I})$ and scaled by a constant $\sigma \in \mathbb{R}_{\geq 0}$. We remark that t represents progress along this interpolation path, not physical time. Notably, stochastic interpolants enable transformations between arbitrary distributions and allow us to assess the performance of the generative process under varying prior distributions q_0 . This contrasts with, e.g., score based diffusion methods [24, 51], which typically assume an isotropic Gaussian prior.

The stochastic interpolant induces a deterministic trajectory of densities $p_t(\mathbf{x})$, governed by an ODE known as the probability flow:

$$d\mathbf{x} = \mathbf{u}_t(\mathbf{x}) dt. \quad (2)$$

If the vector field $\mathbf{u}_t(\mathbf{x})$ was tractable to sample, the weights of a neural network (NN) $\mathbf{v}^\theta(\mathbf{x}, t) : [0, 1] \times \mathbb{R}^d \mapsto \mathbb{R}^d$ could be optimized directly by minimizing

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \mathbf{x} \sim p_t(\mathbf{x})} \left\| \mathbf{u}_t(\mathbf{x}) - \mathbf{v}^\theta(\mathbf{x}, t) \right\|. \quad (3)$$

The learned vector field \mathbf{v}^θ could then be used to generate new samples from the target distribution by starting from $\mathbf{x}_0 \sim q_0$ and integrating the probability flow ODE (Eq. 2), for example, using a numerical scheme such as Euler's method, i.e., $\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \mathbf{v}^\theta(\mathbf{x}_t, t)\Delta t$ for time step Δt .

However, for arbitrary distributions q_0 and q_1 , the objective in Eq. 3 is computationally intractable [52]. Instead, we consider the expectation over interpolated point pairs from the two distributions. Eq. 1 defines a *conditional probability distribution* $p_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathcal{N}(\mathbf{x} | (1-t) \cdot \mathbf{x}_0 + t \cdot \mathbf{x}_1, \sigma^2)$, with *conditional vector field* $\mathbf{u}_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) = \mathbf{x}_1 - \mathbf{x}_0$ [27]. The ability to directly sample from the conditional probability via Eq. 1 allows formulating the conditional flow matching (CFM) objective

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \mathbf{x}_0 \sim q_0, \mathbf{x}_1 \sim q_1, \mathbf{x} \sim p_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1)} \left\| \mathbf{u}_t(\mathbf{x}|\mathbf{x}_0, \mathbf{x}_1) - \mathbf{v}^\theta(\mathbf{x}, t) \right\|^2. \quad (4)$$

Conditioning

4.1 Conditioning

Each DiT block receives conditioning signals that encode information about the time step t , as well as atom- and pair-wise information derived from a molecular graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Pair-wise information is used during the self-attention update (see [Eq. 5](#) and [Sec. 4.3](#)), whereas atom-wise information and the time conditioning signal are injected via adaptive layer norms (see [Eq. 6](#)).

Time conditioning The current time t of the latent state \mathbf{x}_t is encoded via a two-layer MLP as

$$\mathbf{c}^t = \text{MLP}(t). \quad (7)$$

Atom-wise conditioning Atom-wise conditioning tokens are obtained from a GNN inspired by the processor module of the MeshGraphNet (MGN) framework [\[54\]](#)

$$\mathbf{c}_i^{\mathcal{G}} = \text{GNN}_{\text{node}}(\mathcal{V}, \mathcal{E}), \quad (8)$$

where $\mathbf{c}_i^{\mathcal{G}}$ denotes the final node representation for atom i . See [Appendix I](#) for details on the GNN.

Pair-wise conditioning In addition to node representations, the GNN also produces edge-level representations, which can be used to define pair-wise conditioning tokens:

$$\text{bond-pair: } \mathbf{c}_{ij}^{\mathcal{G}} = \begin{cases} \text{GNN}_{\text{edge}}(\mathcal{V}, \mathcal{E}) & \forall (i, j) \in \mathcal{E} \\ \bar{\mathbf{c}}^{\mathcal{G}} & \forall (i, j) \notin \mathcal{E}. \end{cases} \quad (9)$$

These tokens only capture interactions between bonded atoms, i.e., when $(i, j) \in \mathcal{E}$. Conditioning tokens for non-bonded pairs are set to a learnable vector $\bar{\mathbf{c}}^{\mathcal{G}}$. Since self-attention operates on *all* atom pairs (i, j) , even if they are not connected by a chemical bond, we also define an alternative pair-wise conditioning method inspired by the Graphormer architecture [\[55\]](#)

$$\text{all-pair: } \mathbf{c}_{ij}^{\mathcal{G}} = \text{MLP}(s(i, j)), \quad (10)$$

where $s(i, j)$ denotes the graph geodesic (i.e., the shortest path between atoms i and j via edges in \mathcal{G}). This formulation allows conditioning on all atom pairs, even if they are not directly connected. In the following, we refer to these two variants as **bond-pair** ([Eq. 9](#)) and **all-pair** ([Eq. 10](#)) conditioning, respectively (see [Sec. 5](#) for empirical comparisons between both methods).



Positional Embedding

4.2 Positional Embeddings (PE)

We inject information about the current atomic positions $\mathcal{R} = \{\vec{r}_1, \dots, \vec{r}_N \mid \vec{r}_i \in \mathbb{R}^3\}$ of the latent state \mathbf{x}_t via positional embeddings (PEs). To that end, we assume the positions to have zero center of mass (see section A.1). As mentioned in Sec. 3.1, conformer identity is determined only by the *relative* positions of atoms (i.e., it is invariant under rigid translations and rotations), but ensuring that these symmetries are preserved within the network often incurs computational overhead.

Therefore, we examine a representative range of PEs that vary in the number of Euclidean symmetries they respect by construction, which affects how the latent representations transform under translations and rotations. This includes absolute and relative PEs, motivated by a similar distinction made in other application domains [56, 57]. Inspired by recent advances in geometric deep learning [58–60], we also propose a PE strategy, obeying all Euclidean symmetries. Our aim is to find an acceptable tradeoff between model accuracy, efficiency and scalability, by scanning the design space of PE strategies.

Absolute Positional Embeddings (aPE) Following Refs. [31, 37] they are calculated as

$$\mathbf{p}_i^{\text{aPE}} = \text{MLP}(\vec{r}_i), \quad (11)$$

where $\vec{r}_i \in \mathbb{R}^3$ is the Cartesian position of the i -th atom. This kind of PE is neither rotationally nor translationally invariant, and serves as a baseline without any symmetry constraints.

Relative Positional Embedding (rPE) We use displacements vectors $\vec{r}_{ij} = \vec{r}_i - \vec{r}_j$ to build pairwise relative PEs as

$$\mathbf{p}_{ij}^{\text{rPE}} = \text{MLP}(\vec{r}_{ij}). \quad (12)$$

This formulation ensures translational invariance but not rotational invariance.

Euclidean Positional Embedding (PE(3)) Adapting ideas from equivariant message passing neural networks like PaiNN [61] or NequIP [62], we construct SO(3)-equivariant pairwise PEs as a concatenation of $L + 1$ components

$$\mathbf{p}_{ij}^{\text{PE(3)}} = \bigoplus_{\ell=0}^L \phi_{\ell}(r_{ij}) \odot \mathbf{Y}_{\ell}(\hat{r}_{ij}), \quad (13)$$

where $\phi_{\ell} : \mathbb{R} \mapsto \mathbb{R}^{1 \times H}$ is a radial filter function, $\hat{r} = \vec{r}/r$, and $\mathbf{Y}_{\ell} \in \mathbb{R}^{(2\ell+1) \times 1}$ are spherical harmonics of degree $\ell = 0 \dots L$. The element-wise multiplication ‘ \odot ’ between radial filters and spherical harmonics is understood to be “broadcasting” along axes with size 1, such that $(\phi_{\ell} \odot \mathbf{Y}_{\ell}) \in \mathbb{R}^{(2\ell+1) \times H}$ and (after concatenation) $\mathbf{p}_{ij}^{\text{PE(3)}} \in \mathbb{R}^{(L+1)^2 \times H}$. Under rotation of the input positions, these PEs transform equivariantly (see Appendix B.2.2). Moreover, because displacement vectors are used as inputs, the embeddings are also invariant to translations. As a result, they respect the full set of Euclidean symmetries relevant to molecular geometry.

Self-attention

4.3 Self-Attention Operation

For ease of notation, we only describe self-attention with a single head, but employ multi-head attention [29] with n_{heads} heads in our experiments. All self-attention blocks rely on query, key and value vectors, which are obtained from the input tokens $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_N \mid \mathbf{h}_i \in \mathbb{R}^H\}$ as

$$\mathbf{q} = \mathbf{W}_q \tilde{\mathbf{h}}, \quad \mathbf{k} = \mathbf{W}_k \tilde{\mathbf{h}}, \quad \mathbf{v} = \mathbf{W}_v \tilde{\mathbf{h}}, \quad (14)$$

where $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v \in \mathbb{R}^{H \times H}$ are trainable weight matrices and $\tilde{\mathbf{h}}$ is either identical to \mathbf{h} or combines it with a PE (see below). We define a slightly modified similarity kernel

$$\text{sim}(\mathbf{q}, \mathbf{k}, \mathbf{u}) = \exp \left(\frac{\mathbf{q}^\top \cdot (\mathbf{k} \odot \mathbf{u})}{\sqrt{H}} \right), \quad (15)$$

where $\mathbf{u} \in \mathbb{R}^H$ is used to inject additional information, e.g., conditioning signals and/or positional embeddings, and ‘ \odot ’ denotes element-wise multiplication. Depending on the subset of Euclidean symmetries we aim to incorporate, we adopt slightly different formulations of the self-attention mechanism, as detailed below.

Standard Self-Attention For absolute and relative PEs, we slightly modify standard self-attention to allow injecting pair-wise information into the values in addition to using our modified similarity kernel:

$$\text{ATT}(\mathcal{H})_i = \frac{\sum_{j=1}^N \text{sim}(\mathbf{q}_i, \mathbf{k}_j, \mathbf{u}_{ij}) \cdot (\mathbf{v}_j \odot \mathbf{u}_{ij})}{\sum_{j=1}^N \text{sim}(\mathbf{q}_i, \mathbf{k}_j, \mathbf{u}_{ij})}. \quad (16)$$

Queries, keys, and values are obtained with [Eq. 14](#) from different (position-encoded) tokens $\tilde{\mathbf{h}}_i$ depending on the chosen PEs; further, the injected pair-wise information \mathbf{u}_{ij} differs:

$$\tilde{\mathbf{h}}_i = \begin{cases} \mathbf{h}_i + \mathbf{p}_i^{\text{aPE}} & \text{for absolute PEs,} \\ \mathbf{h}_i & \text{for relative PEs,} \end{cases} \quad \text{and} \quad \mathbf{u}_{ij} = \begin{cases} \mathbf{c}_{ij}^{\mathcal{G}} & \text{for absolute PEs,} \\ \mathbf{c}_{ij}^{\mathcal{G}} + \mathbf{p}_{ij}^{\text{rPE}} & \text{for relative PEs.} \end{cases} \quad (17)$$

Here $\mathbf{c}_{ij}^{\mathcal{G}} \in \mathbb{R}^H$ are pair-wise graph conditioning tokens (see [Eqs. 9](#) and [10](#)) and $\mathbf{p}_i^{\text{aPE}} \in \mathbb{R}^H$ and $\mathbf{p}_{ij}^{\text{rPE}} \in \mathbb{R}^H$ are the absolute and relative PEs described above (see [Eqs. 11](#) and [12](#)).

SO(3)-Equivariant Self-Attention To preserve all Euclidean symmetries throughout the network, every token must transform equivariantly. One way to achieve this is by separating out the rotational degrees of freedom, encoding them with irreducible representations of the rotation group $\text{SO}(3)$. This introduces a “degree-axis” of size $(L+1)^2$, which encodes angular components of increasing order. The maximum degree L is chosen to ensure high fidelity at a reasonable computational cost. For example, setting $L = 1$ restricts the representation to scalars and vectors, as used in models like PaiNN [\[61\]](#) or TorchMDNet [\[34\]](#). An $\text{SO}(3)$ -equivariant formulation of self-attention is then given as

$$\text{ATT}_{\text{SO}(3)}(\mathcal{H})_i = \frac{\sum_{j=1}^N \text{sim}(\mathbf{q}_i, \mathbf{k}_j, \mathbf{u}_{ij}) \cdot (\hat{\mathbf{u}}_{ij} \otimes \mathbf{v}_j)}{\sum_{j=1}^N \text{sim}(\mathbf{q}_i, \mathbf{k}_j, \mathbf{u}_{ij})}, \quad (18)$$

where equivariant queries, keys and values can be calculated similarly to [Eq. 14](#) and ‘ \otimes ’ denotes a Clebsch-Gordan (CG) tensor product contraction [\[60\]](#). The dot-product in the similarity measure is taken along both feature and degree axes, such that the overall update preserves equivariance (see Appendix for details). Tokens and scaling vectors are calculated as

$$\tilde{\mathbf{h}}_i = \mathbf{h}_i, \quad \mathbf{u}_{ij} = \phi(r_{ij}) \odot \mathbf{c}_{ij}^{\mathcal{G}}, \quad \hat{\mathbf{u}}_{ij} = \mathbf{p}_{ij}^{\text{PE}(3)} \odot \mathbf{c}_{ij}^{\mathcal{G}}, \quad (19)$$

where $\phi(r_{ij}) \in \mathbb{R}^{(L+1)^2 \times H}$ is a radial filter, and the element-wise products with the pair-wise conditioning tokens $\mathbf{c}_{ij}^{\mathcal{G}} \in \mathbb{R}^{1 \times H}$ are broadcast along the degree axis. Importantly, the $2\ell + 1$ subcomponents of the radial filter for degree ℓ are obtained by repeating per-degree filter functions $\phi_{\ell}(r_{ij}) \in \mathbb{R}^{1 \times H}$ along the degree axis to preserve equivariance (see also [Eq. 13](#)).

Thank You