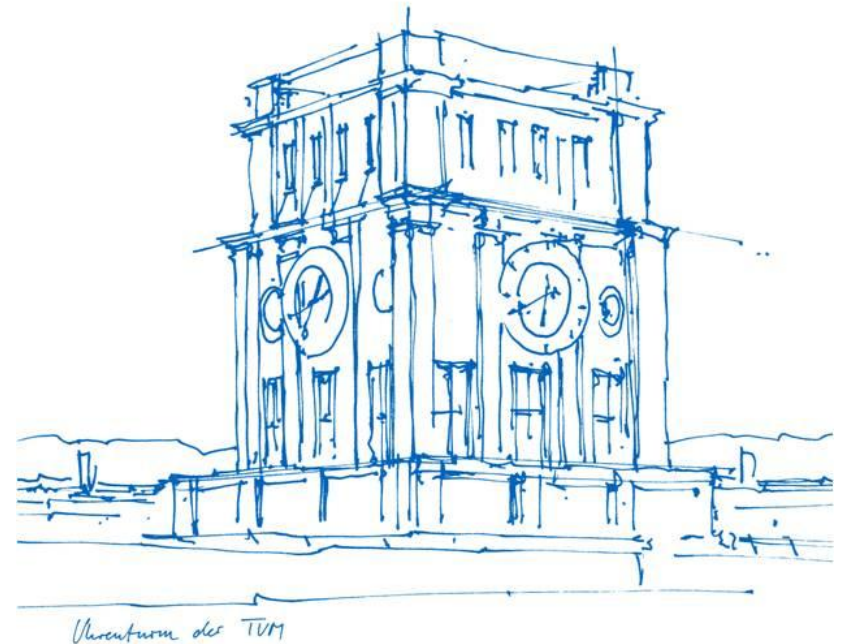


Equivariant Spatio-Temporal Attentive Graph Networks for Physical Dynamics

Celia Tundidor Centeno



Previously...

What are we discussing today?

How can we use machine learning to simulate physical systems with **high fidelity to their dynamics**?

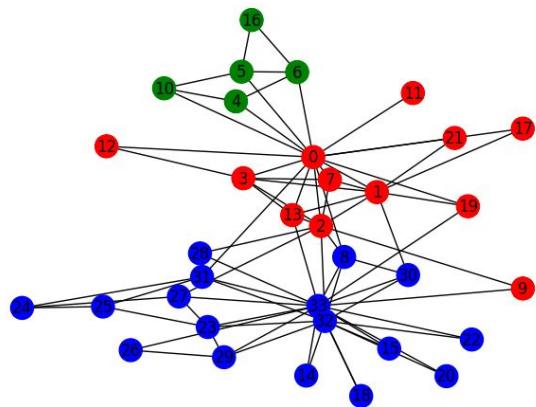
Contexts of application:

- topics: molecular dynamics, protein structure prediction, robotics...
- levels: macro, protein, smaller molecules...

Foundational concepts: Graph Neural Networks (GNNs)

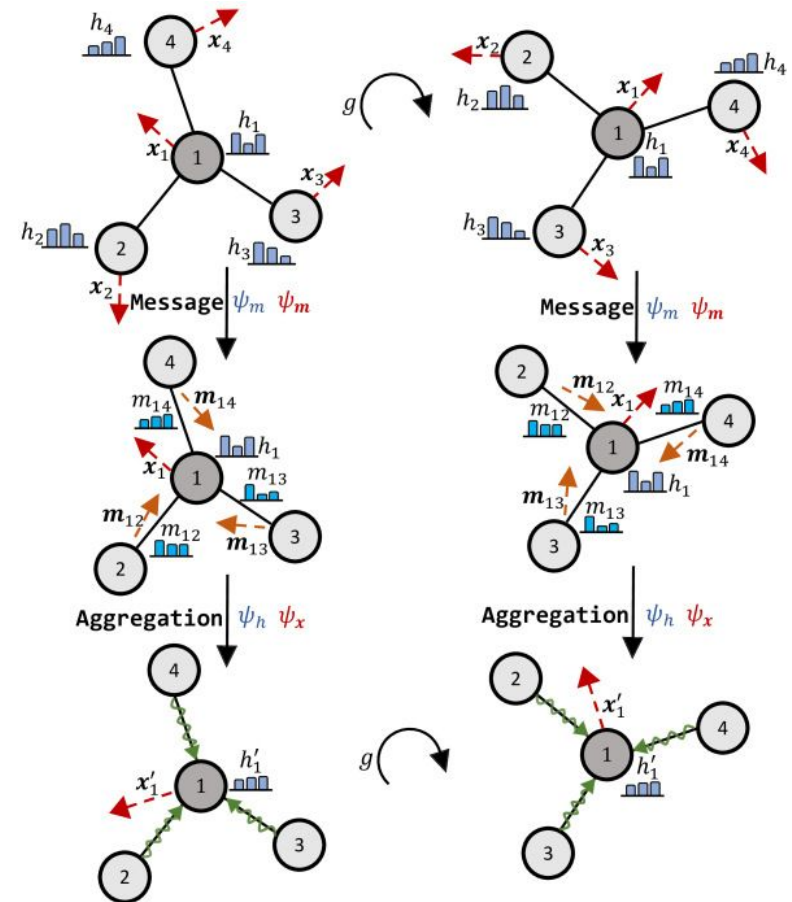
Naturally fit for physical system representation

- Unit elements as nodes (e.g., atoms)
- Relations as edges (e.g., chemical bonds)
- Latent interactions as message passing between these nodes with edges



Foundational concepts: Equivariance

Output reflects a predictable transformation equivalent to that of the input. Physical consistence irrespective of the coordinate system and view



State of the Art: equivariant GNNs

Spatially: generalising GNNs to fit the symmetry of our world

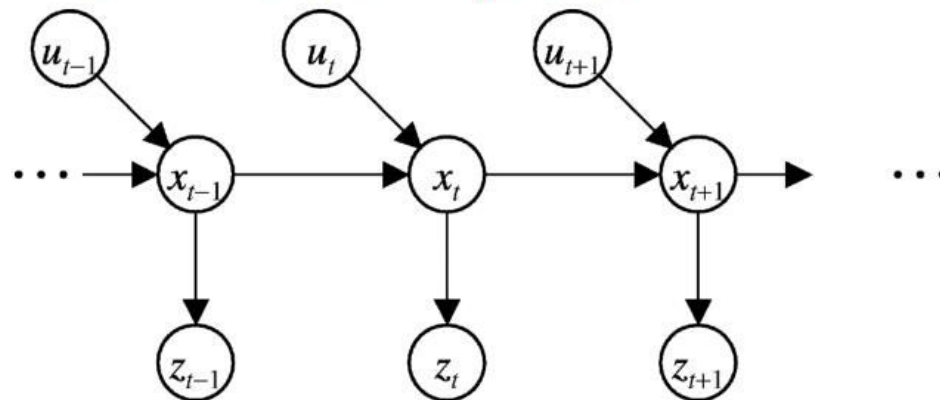
Temporally: **frame-to-frame forecasting**

E.g.:

- Tensor-Field Networks (TFN)
- SE(3)-Transformer
- LieTransformer and LieConv
- E(n)-equivariant GNNs (EGNN)
- Equivariant Graph Mechanics Networks (GMN)

The problem: the Markovian assumption

“The future state only depends on the current state, independent of all other past states”



$$p(z_t | x_{0:t}, z_{1:t}, u_{1:t}) = p(z_t | x_t)$$

$$p(x_t | x_{1:t-1}, z_{1:t}, u_{1:t}) = p(x_t | x_{t-1}, u_t)$$

The problem: the Markovian assumption

“The future state only depends on the current state, independent of all other past states”

Previous methods rely on this:

- A single input: system's conformation at a single frame.
- A fixed time step: they predict the future after a fixed time interval (frame-to-frame)

Why is the Markovian assumption problematic?

What if there are unobserved objects interacting with the system?

- Missed by the last frame
- Untracked

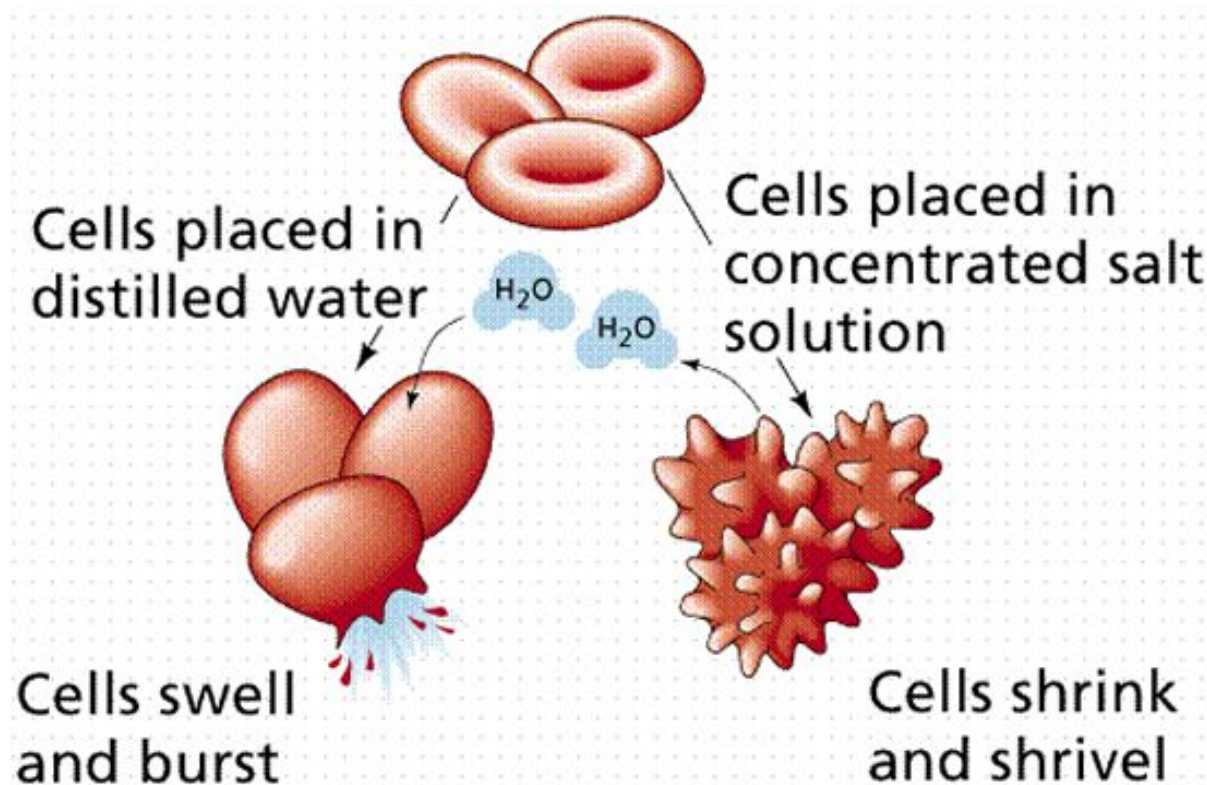
What if the effects induced by other objects are not constant or linear?



Why is the Markovian assumption problematic?

For molecular dynamics in particular:

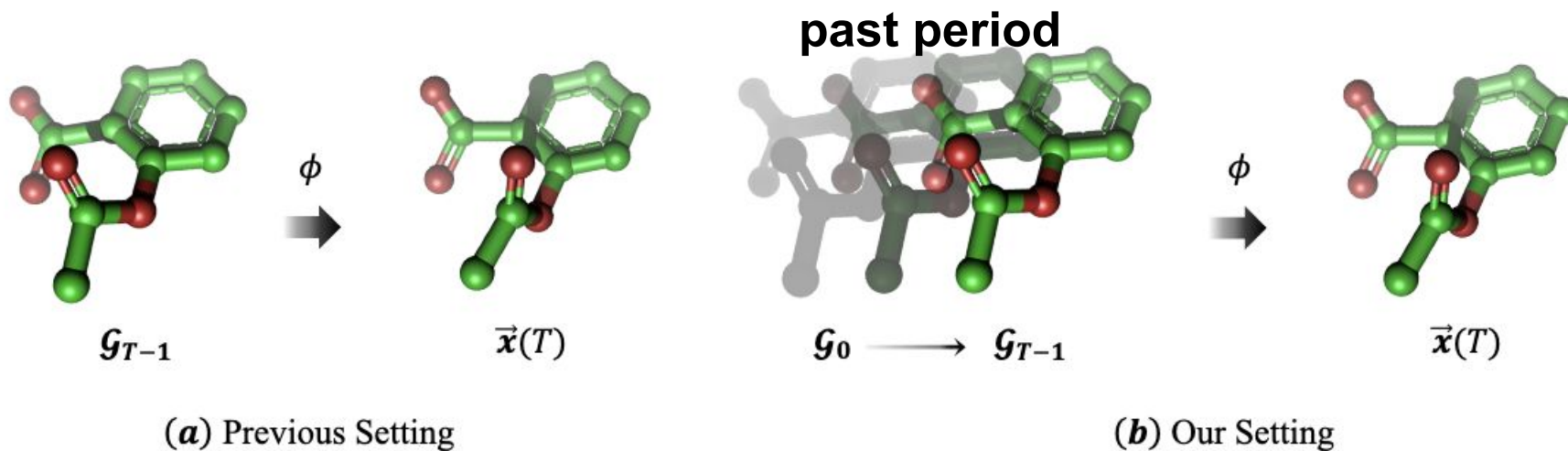
What about solvents (untracked object)?



Addressing Non-Markovian Dynamics

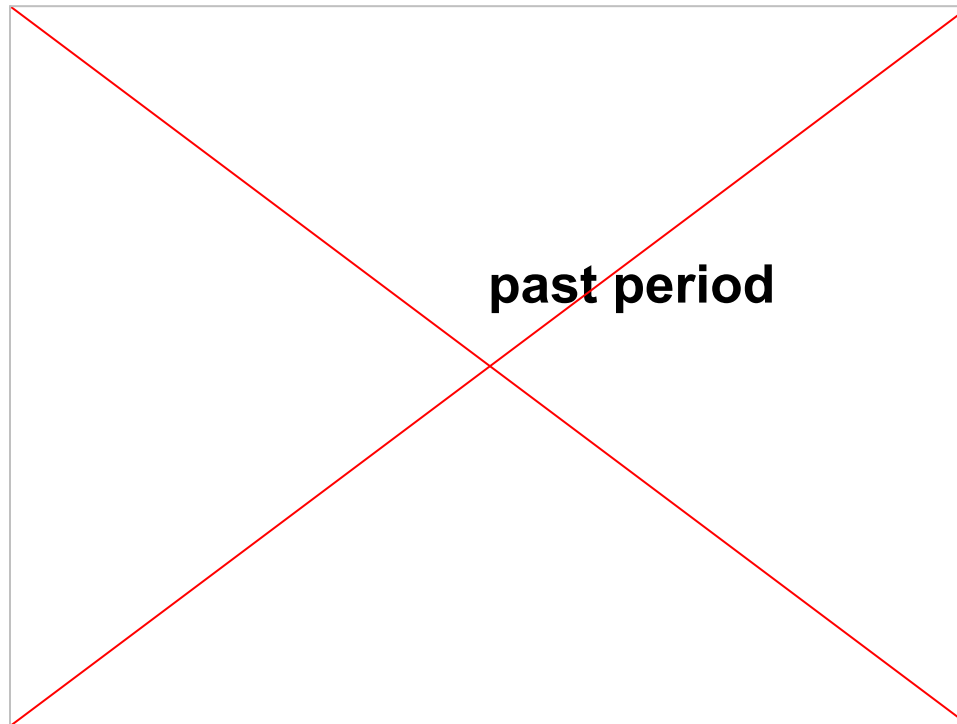
We define a **past period** (to be taking as input)

Idea: if the past period is sufficiently long, non-Markovian behaviour can be recovered



Addressing Non-Markovian Dynamics

We can also recover periodic motion (e.g. periodic thermal vibration)



Addressing Non-Markovian Dynamics

We can therefore use Spatio-Temporal Graph Neural Networks

(STGNNs)...but they are unfit for Euclidean symmetry and physical laws

- traditional use case not on physical modelling (e.g. traffic forecasting)
- no 3D geometric equivariance

Enter ESTAG

Equivariant Spatio-Temporal Attentive Graph Networks (ESTAG):

- capturing non-Markovian behaviour (based on STGNNs)
- making STGNNs equivariant (for Euclidean symmetry)

ESTAG components

1. Equivariant Discrete Fourier Transform (EDFT): extracts periodic patterns
2. Equivariant Spatial Module (ESM): passes spatial messages.
3. Equivariant Temporal Module (ETM): aggregates temporal messages using forward attention and equivariant pooling

Equivariant Discrete Fourier Transform (EDFT)

Fourier Transform helps us understand the frequency domain

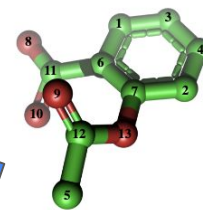
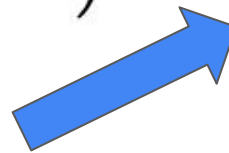
-> periodicity (node-wise temporal dynamics for the global context). c_i is the frequency amplitude of node i .

We can later use this information to check node cross-correlation (A)

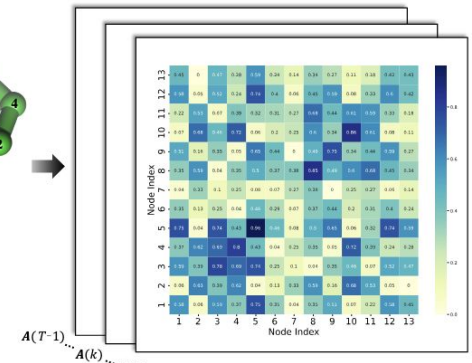
A and c are $E(3)$ -invariant!

$$\vec{f}_i(k) = \sum_{t=0}^{T-1} e^{-i' \frac{2\pi}{T} kt} \left(\vec{x}_i(t) - \overline{\vec{x}(t)} \right)$$

$$c_i(k) = w_k(\mathbf{h}_i) \|\vec{f}_i(k)\|^2.$$



Aspirin



$$\mathbf{A}_{ij}(k) = w_k(\mathbf{h}_i)w_k(\mathbf{h}_j)|\langle \vec{f}_i(k), \vec{f}_j(k) \rangle|$$

Equivariant Spatial Module (ESM)

Encoding and passing the spatial geometry of each graph through each layer

EGNN + EDFT features:

- + correlation (A_{ij}) to evaluate global temporal connections
- + amplitude (c_i) to update hidden features at each node

Equivariant Spatial Module (ESM)

Process: compute messages, update hidden features, update positions

$$\mathbf{m}_{ij} = \phi_m \left(\mathbf{h}_i^{(l)}(t), \mathbf{h}_j^{(l)}(t), \|\vec{\mathbf{x}}_{ij}^{(l)}(t)\|^2, \mathbf{A}_{ij} \right),$$

$$\mathbf{h}_i^{(l+1)}(t) = \mathbf{h}_i^{(l)}(t) + \phi_h \left(\mathbf{h}_i^{(l)}(t), \mathbf{c}_i, \sum_{j \neq i} \mathbf{m}_{ij} \right),$$

$$\vec{\mathbf{a}}_i(t) = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \vec{\mathbf{x}}_{ij}^{(l)}(t) \phi_x(\mathbf{m}_{ij}),$$

$$\vec{\mathbf{x}}_i^{(l+1)}(t) = \vec{\mathbf{x}}_i^{(l)}(t) + \vec{\mathbf{a}}_i(t),$$

note that these operations do not disturb equivariance

Equivariant Temporal Module (ETM)

Modelling self-correspondence with an attention mechanism

Forward temporal attention: we only rely on the past

Equivariant pooling

Equivariant Temporal Module (ETM)

Modelling self-correspondence with an attention mechanism

Forward temporal attention: we only rely on the past

Equivariant pooling: aggregates spatial and temporal information

$$\alpha_i^{(l)}(ts) = \frac{\exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))}{\sum_{s=0}^t \exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))}, \quad \text{attention weight}$$

$$\mathbf{h}_i^{(l+1)}(t) = \mathbf{h}_i^{(l)}(t) + \sum_{s=0}^t \alpha_i^{(l)}(ts) \mathbf{v}_i^{(l)}(s), \quad \text{hidden feature}$$

temporal displacement vector

$$\vec{\mathbf{x}}_i^{(l+1)}(t) = \vec{\mathbf{x}}_i^{(l)}(t) + \sum_{s=0}^t \alpha_i^{(l)}(ts) \vec{\mathbf{x}}_i^{(l)}(ts) \phi_x(\mathbf{v}_i^{(l)}(s)),$$

Equivariant Temporal Pooling

Equivariant pooling: apply a linear transformation to the updated coordinates

$$\vec{x}_i^*(T) = \hat{\mathbf{X}}_i \mathbf{w} + \vec{x}_i^{(L)}(T-1),$$



Training via MSE loss

$$\mathcal{L} = \sum_{i=1}^N \|\vec{x}_i(T) - \vec{x}_i^*(T)\|_2^2.$$

Architecture recap

Input: historical series of spatio-temporal graphs $\{G_t\}$ from time $t=0$ to $T-1$

Equivariant Discrete Fourier Transform (EDFT): processes historical trajectory for each node. Extracts equivariant frequency features \rightarrow invariant node features (c) and adjacency matrix (A).

Stacked Modules: computes spatial and temporal relationships. L layers of alternating equivariant components (ESM, ETM)

Equivariant Temporal Pooling: pooling layer to combine time and space dependencies

Output: position of each node at time T

Architecture recap

EDFT:

$$\vec{f}_i(k) = \sum_{t=0}^{T-1} e^{-i' \frac{2\pi}{T} kt} \left(\vec{x}_i^\alpha(t) - \overline{\vec{x}^\alpha(t)} \right),$$

$$\mathbf{A}_{ij}(k) = w_k(\mathbf{h}_i) w_k(\mathbf{h}_j) |\langle \vec{f}_i(k), \vec{f}_j(k) \rangle|,$$

$$\mathbf{c}_i(k) = w_k(\mathbf{h}_i) \|\vec{f}_i(k)\|^2.$$

Architecture recap

ESM:

$$\begin{aligned}\mathbf{m}_{ij} &= \phi_m \left(\mathbf{h}_i^{(l)}(t), \mathbf{h}_j^{(l)}(t), \frac{(\vec{\mathbf{X}}_{ij}^{(l)}(t))^\top \vec{\mathbf{X}}_{ij}^{(l)}(t)}{\|(\vec{\mathbf{X}}_{ij}^{(l)}(t))^\top \vec{\mathbf{X}}_{ij}^{(l)}(t)\|_F}, \mathbf{A}_{ij} \right), \\ \mathbf{h}_i^{(l+1)}(t) &= \phi_h \left(\mathbf{h}_i^{(l)}(t), \mathbf{c}_i(k), \sum_{j \neq i} \mathbf{m}_{ij} \right), \\ \vec{\mathbf{A}}_i^{(l)}(t) &= \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \vec{\mathbf{X}}_{ij}^{(l)}(t) \phi_{\mathbf{X}}(\mathbf{m}_{ij}), \\ \vec{\mathbf{X}}_i^{(l+1)}(t) &= \vec{\mathbf{X}}_i^{(l)}(t) + \vec{\mathbf{A}}_i^{(l)}(t).\end{aligned}$$

ETM:

$$\begin{aligned}\alpha_i^{(l)}(ts) &= \frac{\exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))}{\sum_{s=0}^t \exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))}, \\ \mathbf{h}_i^{(l+1)}(t) &= \mathbf{h}_i^{(l)}(t) + \sum_{s=0}^t \alpha_i^{(l)}(ts) \mathbf{v}_i^{(l)}(s), \\ \vec{\mathbf{X}}_i^{(l+1)}(t) &= \vec{\mathbf{X}}_i^{(l)}(t) + \sum_{s=0}^t \alpha_i^{(l)}(ts) \vec{\mathbf{X}}_i^{(l)}(ts) \phi_{\mathbf{X}}(\mathbf{v}_i^{(l)}(s)),\end{aligned}$$

where

$$\begin{aligned}\mathbf{q}_i^{(l)}(t) &= \phi_q \left(\mathbf{h}_i^{(l)}(t) \right), \\ \mathbf{k}_i^{(l)}(t) &= \phi_k \left(\mathbf{h}_i^{(l)}(t) \right), \\ \mathbf{v}_i^{(l)}(t) &= \phi_v \left(\mathbf{h}_i^{(l)}(t) \right).\end{aligned}$$

Equivariance details

Theorem A.1. We denote ESTAG as $\vec{X}(T) = \phi \left(\{(\mathbf{H}(t), g \cdot \vec{X}(t), \mathbf{A})\}_{t=0}^{T-1} \right)$, then ϕ is E(3)-equivariant.

Proof. 1. We firstly prove that EDFT is E(3)-equivariant.

$$\mathbf{O}\vec{f}_i(k) = \sum_{t=0}^{T-1} e^{-i' \frac{2\pi}{T} kt} \left(\mathbf{O}\vec{x}_i(t) + \mathbf{b} - \overline{\mathbf{O}\vec{x}(t) + \mathbf{b}} \right),$$

$$\mathbf{A}_{ij}(k) = w_k(\mathbf{h}_i)w_k(\mathbf{h}_j)|\langle \mathbf{O}\vec{f}_i(k), \mathbf{O}\vec{f}_j(k) \rangle|,$$

$$\mathbf{c}_i(k) = w_k(\mathbf{h}_i)\|\mathbf{O}\vec{f}_i(k)\|^2.$$

2. We secondly prove the E(3)-equivariance of ESM.

$$\mathbf{m}_{ij} = \phi_m \left(\mathbf{h}_i^{(l)}(t), \mathbf{h}_j^{(l)}(t), \|\mathbf{O}\vec{x}_{ij}^{(l)}(t)\|^2, \mathbf{A}_{ij} \right),$$

$$\mathbf{h}_i^{(l+1)}(t) = \phi_h \left(\mathbf{h}_i^{(l)}(t), \mathbf{c}_i(k), \sum_{j \neq i} \mathbf{m}_{ij} \right),$$

$$\mathbf{O}\vec{a}_i(t) = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \mathbf{O}\vec{x}_{ij}^{(l)}(t) \phi_x(\mathbf{m}_{ij}),$$

$$\mathbf{O}\vec{x}_i^{(l+1)}(t) + \mathbf{b} = \mathbf{O}\vec{x}_i^{(l)}(t) + \mathbf{b} + \mathbf{O}\vec{a}_i^{(l+1)}(t).$$

Equivariance details

3. We then prove that ETM is E(3)-equivariant.

$$\mathbf{q}_i^{(l)}(t) = \phi_q \left(\mathbf{h}_i^{(l)}(t) \right),$$

$$\mathbf{k}_i^{(l)}(t) = \phi_k \left(\mathbf{h}_i^{(l)}(t) \right),$$

$$\mathbf{v}_i^{(l)}(t) = \phi_v \left(\mathbf{h}_i^{(l)}(t) \right),$$

$$\alpha_i^{(l)}(ts) = \frac{\exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))}{\sum_{s=0}^t \exp(\mathbf{q}_i^{(l)}(t)^\top \mathbf{k}_i^{(l)}(s))},$$

$$\mathbf{h}_i^{(l+1)}(t) = \mathbf{h}_i^{(l)}(t) + \sum_{s=0}^t \alpha_i^{(l)}(ts) \mathbf{v}_i^{(l)}(s),$$

$$\mathbf{O}\vec{\mathbf{x}}_i^{(l+1)}(t) + \mathbf{b} = \mathbf{O}\vec{\mathbf{x}}_i^{(l)}(t) + \mathbf{b} + \sum_{s=0}^t \alpha_i^{(l)}(ts) \mathbf{O}\vec{\mathbf{x}}_i^{(l)}(ts) \phi_x(\mathbf{v}_i^{(l)}(s)).$$

4. We finally prove that the linear pooling is equivariant:

$$\mathbf{O}\vec{\mathbf{x}}_i^*(T) + \mathbf{b} = \mathbf{O}\hat{\mathbf{X}}_i \mathbf{w} + \mathbf{O}\vec{\mathbf{x}}_i^{(L)}(T-1) + \mathbf{b}.$$

Experiments

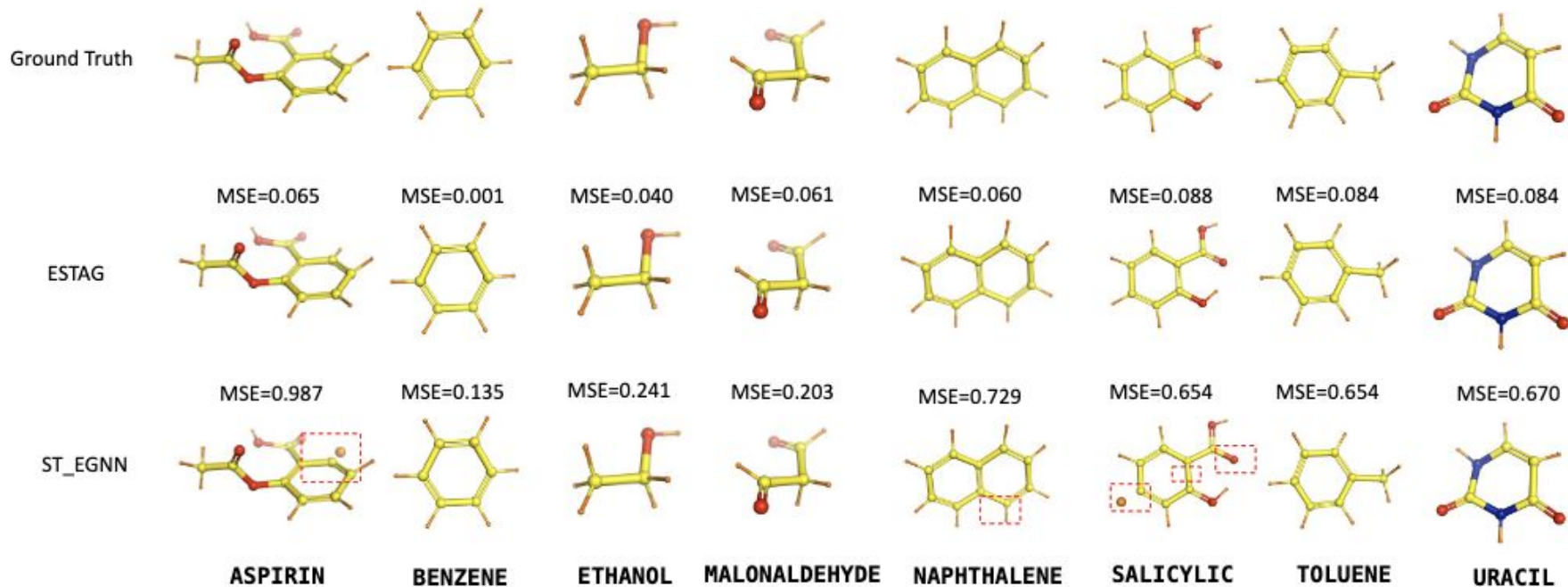
Testing on three datasets for the different levels:

Molecular: MD17, trajectories of small molecules (e.g., Aspirin, Benzene, Ethanol) generated by Molecular Dynamics simulation. External temperature and pressure are unobserved (non-Markovian behaviour)

Protein-level: AdK equilibrium trajectory dataset (protein dynamics). The dynamics of water and ions are unobserved (non-Markovian behaviour)

Macro-level: CMU Motion Capture Database (human motion trajectories) (e.g., walking, basketball). Environmental states are unobserved (non-Markovian behaviour)

Experimental results: molecular



Experimental results: molecular

Table 1: Prediction error ($\times 10^{-3}$) on MD17 dataset. Results averaged across 3 runs. We do not display the standard deviation due to its small value.

	ASPIRIN	BENZENE	ETHANOL	MALONALDEHYDE	NAPHTHALENE	SALICYLIC	TOLUENE	URACIL
PT- <i>s</i>	15.579	4.457	4.332	13.206	8.958	12.256	6.818	10.269
PT- <i>m</i>	9.058	2.536	2.688	6.749	6.918	8.122	5.622	7.257
PT- <i>t</i>	0.715	0.114	0.456	0.596	0.737	0.688	0.688	0.674
EGNN- <i>s</i>	12.056	3.290	2.354	10.635	4.871	8.733	3.154	6.815
EGNN- <i>m</i>	6.237	1.882	1.532	4.842	3.791	4.623	2.516	3.606
EGNN- <i>t</i>	0.625	0.112	0.416	0.513	0.614	0.598	0.577	0.568
ST_TFN	0.719	0.122	0.432	0.569	0.688	0.684	0.628	0.669
ST_GNN	1.014	0.210	0.487	0.664	0.769	0.789	0.713	0.680
ST_SE(3) _{TR}	<u>0.669</u>	0.119	0.428	0.550	0.625	0.630	0.591	0.597
ST_EGNN	<u>0.735</u>	0.163	<u>0.245</u>	<u>0.427</u>	0.745	0.687	<u>0.553</u>	<u>0.445</u>
EqMOTION	0.721	0.156	0.476	0.600	0.747	0.697	0.691	0.681
STGCN	0.715	<u>0.106</u>	0.456	0.596	0.736	0.682	0.687	0.673
AGL-STAN	0.719	<u>0.106</u>	0.459	0.596	<u>0.601</u>	<u>0.452</u>	0.683	0.515
ESTAG	0.063	0.003	0.099	0.101	0.068	0.047	0.079	0.066

Experimental results: protein and macro

METHOD	MSE	TIME(S)
PT- <i>s</i>	3.260	-
PT- <i>m</i>	3.302	-
PT- <i>t</i>	2.022	-
EGNN- <i>s</i>	3.254	1.062
EGNN- <i>m</i>	3.278	1.088
EGNN- <i>t</i>	1.983	1.069
ST_GNN	1.871	2.769
ST_GMN	<u>1.526</u>	4.705
ST_EGNN	1.543	4.705
STGCN	1.578	1.840
AGL-STAN	1.671	1.478
ESTAG	1.471	6.876

METHOD	WALK	BASKETBALL
PT- <i>s</i>	329.474	886.023
PT- <i>m</i>	127.152	413.306
PT- <i>t</i>	3.831	15.878
EGNN- <i>s</i>	63.540	749.486
EGNN- <i>m</i>	32.016	335.002
EGNN- <i>t</i>	0.786	12.492
ST_GNN	0.441	15.336
ST_TFN	0.597	13.709
ST_SE(3)TR	0.236	13.851
ST_EGNN	0.538	13.199
EqMOTION	1.011	<u>4.893</u>
STGCN	0.062	4.919
AGL-STAN	0.037	5.734
ESTAG	<u>0.040</u>	0.746

Table 4: Ablation studies ($\times 10^{-3}$) on MD17 dataset. Results averaged across 3 runs.

	Aspirin	Benzene	Ethanol	Malonaldehyde	Naphthalene	Salicylic	Toluene	Uracil
ESTAG	0.063	0.003	0.099	0.101	0.068	0.047	0.079	0.066
w/o EDFT	0.079	0.003	0.108	0.148	0.104	0.145	0.102	0.063
w/o Attention	0.087	0.004	0.104	0.112	0.129	0.095	0.097	0.078
w/o Equivariance	0.762	0.114	0.458	0.604	0.738	0.698	0.690	0.680
w/o Temporal	0.084	0.003	0.111	0.139	0.141	0.098	0.153	0.071

Table 5: MSE on Ethanol *w.r.t.* the number of layers L .

L	1	2	3	4	5	6
MSE ($\times 10^{-4}$)	1.25	0.990	1.096	1.022	1.042	1.028

Ablation studies

Without EDFT: considerably worse performance. wk (learnable) shown to be beneficial as a spectral filter

Without attention: slightly worse performance

Without equivariance: considerably worse performance

Without temporal pooling: slightly worse performance

Paper analysis: contributions and advantages

- Time: modelling non-Markovian features, capturing periodicity, via EDFT and attention mechanism
- Space: Euclidean symmetry
- Good overall performance

Paper analysis: limitations and criticism

- Limited equivariance: missing embedded physical laws, e.g. no conservation of energy
- Limited benchmarks
- Inconsistent baseline comparisons (due to modifications)
- Ablation study interpretations (limited time effects?)
- Visualization as cherry-picking?

Thanks for listening!

Questions?