
Permutation Equivariant Neural Controlled Differential Equations for Dynamic Graph Representation Learning

Torben Berndt¹ Benjamin Walker² Tiexin Qin³
Jan Stühmer^{1,4} Andrey Kormilitzin⁵

¹Heidelberg Institute for Theoretical Studies, Heidelberg, Germany

²Mathematical Institute, University of Oxford, UK

³City University of Hong Kong, Hong Kong

⁴IAR, Karlsruhe Institute of Technology, Karlsruhe, Germany

⁵Department of Psychiatry, Warneford Hospital, Oxford, UK

Abstract

Dynamic graphs exhibit complex temporal dynamics due to the interplay between evolving node features and changing network structures. Recently, Graph Neural Controlled Differential Equations (Graph Neural CDEs) successfully adapted Neural CDEs from paths on Euclidean domains to paths on graph domains. Building on this foundation, we introduce *Permutation Equivariant Neural Graph CDEs*, which project Graph Neural CDEs onto permutation equivariant function spaces. This significantly reduces the model’s parameter count without compromising representational power, resulting in more efficient training and improved generalisation. We empirically demonstrate the advantages of our approach through experiments on simulated dynamical systems and real-world tasks, showing improved performance in both interpolation and extrapolation scenarios.

1 Introduction

Graph Neural Networks (GNNs) [56, 55, 38, 62] have emerged as a leading framework for modelling graph-structured data, demonstrating significant success in applications such as protein folding [31], social recommender systems [22] or traffic forecasting [30]. However, real-world graphs are often dynamic. Protein-protein interactions vary over time due to cellular processes, social networks evolve as relationships shift, and roads close due to building works and car crashes. Effectively capturing these temporal dynamics is crucial for accurate modelling and robust predictions.

For the last century, differential equations have been the cornerstone of modelling continuous change. However, in recent years, deep learning has revolutionised data analysis with its ability to learn complex patterns from vast amounts of data. While these approaches initially developed in parallel, the concept of Neural Differential Equations (NDEs) [46, 51, 7, 34] has bridged the gap, demonstrating their interconnectedness.

Building upon the approach introduced in [48], we propose *Permutation Equivariant Neural Graph Controlled Differential Equations* (PENG-CDEs), a novel framework which addresses key limitations of prior models.¹ Our primary contributions are as follows:

¹The source code is available at: https://anonymous.4open.science/r/perm_equiv_gn_cdes-BBF8/README.md

- Motivated by temporal and spatial symmetries, we derive PENG-CDEs from first principles. Our framework strikes a balance between expressiveness and parameter efficiency (Section 3).
- We formalise this trade-off and prove that, under simplified assumptions, our proposed model is the optimal approximation to Graph Neural CDEs within the space of permutation-equivariant functions (Section 3.2, Theorem 3.1).
- We further prove that the resulting models are equivariant under both time reparametrisations and permutations of the node set (Section 3.2, Proposition 3.2).
- We empirically validate the effectiveness of PENG-CDEs on dynamic node- and graph-level tasks, consistently outperforming both differential equation-based and other spatio-temporal baselines. In particular, our model sets a new state-of-the-art on the TGB-genre node affinity prediction task (Section 4).

1.1 Related work

Continuous-depth GNNs. Recent work has linked GNNs to continuous-time dynamics on graphs by interpreting hidden layers as a discretised time axis. [6, 53] study convolutional GNNs as discretisations of heat diffusion processes on graphs, while [18] proposes architectures that incorporate both parabolic and hyperbolic PDE terms. Reaction-diffusion dynamics are explored in [12, 17], and advection terms are added in [19]. Higher-order temporal derivatives are considered in [20, 21]. [47] proposes a graph-based formulation of Neural ODEs [7], and Neural SDEs are extended to graph domains in [2]. While these models typically assume static graphs, some work addresses dynamic settings. [11, 10] and [48] extend Neural CDEs to graph domains, by modelling node-level time series and graph structure hierarchically, and by treating the evolving topology as a control, respectively.

Spatio-temporal GNNs. A well-established line of work in dynamic graph representation learning models spatial and temporal dynamics using separate modules for each. DCRNN [40] combines a diffusion-based GCN [38] with a recurrent GRU [9] in an encoder-decoder architecture. STGCN [66] applies GCN [16] for spatial modelling and 1D CNNs for temporal processing. Attention-based extensions such as ASTGCN [26] and ASTGNN [27] introduce distinct spatial and temporal attention mechanisms. WaveNet-GCN [42] fuses graph and temporal convolutions, while STIDGCN [41] uses interactive learning to handle heterogeneous time scales.

Temporal Graph Networks. Another approach focuses on event-based message passing over time-stamped interactions. TGN [52] introduced this framework, later extended by TGNv2 [58] which considers an identification between source and target nodes. Node-centric methods such as DyRep [59] and TCL [64] compute embeddings from temporal neighbourhoods and aggregate across edges. In contrast, edge-centric methods like CAWN [65] and GraphMixer [14] embed edge interactions directly for prediction.

2 Background

2.1 (Temporal) Graph representation learning

We denote the set $\{1, \dots, n\}$ by $[n]$ and the set of functions from X to Y by $\chi(X, Y)$.

Graphs. A *graph* $G = (\mathcal{V}, \mathcal{E})$ consists of a collection of *nodes* $\mathcal{V} = [n]$ and *edges* $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. The graph topology is described by an *adjacency matrix* $\mathbf{A} \in \mathbb{R}^{n \times n}$, where $\mathbf{A}^{ij} = 1$ if $(i, j) \in \mathcal{E}$ and 0 otherwise. The *degree* $d_i = \sum_{j \in \mathcal{V}} \mathbf{A}^{ij}$ of a node is the number of edges incident on it. Let \mathbf{D} be the diagonal *degree matrix* with $\mathbf{D}^{i,i} = d_i$. The *graph Laplacian* is defined as $\mathbf{L} = \mathbf{D} - \mathbf{A}$ and the *normalised graph Laplacian* is defined as $\mathcal{L} = \mathbf{I}_n - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ where \mathbf{I}_n is the $n \times n$ identity matrix. We assume the nodes are equipped with d_x -dimensional *node features* $\mathbf{x}_i \in \mathbb{R}^{d_x}$, stacked to form a node attribute matrix $\mathbf{X} \in \mathbb{R}^{n \times d_x}$. The goal of *Graph Representation Learning* is to learn a latent node representation $\mathbf{Z} \in \mathbb{R}^{n \times d_z}$ which embeds the nodes into some Euclidean space.

Temporal Graph Representation Learning. In this setting, we observe a sequence of graph snapshots $\mathcal{G} = ((t_0, G_{t_0}), \dots, (t_N, G_{t_N}))$ generated by some unknown continuous-time underlying process. Each snapshot (t_k, G_{t_k}) captures the graph state at time t_k , with $G_{t_k} = (\mathcal{V}, \mathcal{E}_{t_k})$ and

corresponding adjacency matrices $\mathbf{A}_{t_k} \in \mathbb{R}^{n \times n}$. The objective is to learn a dynamic non-linear latent representation $\mathbf{Z}_t \in \mathbb{R}^{n \times d_z}$ for the nodes for all times $t \in (t_0, t_N]$.

Equivariance. Let G be a group and X, Y be two sets with a (left) G -action. A function $f : X \rightarrow Y$ is G -equivariant if $f(g \cdot x) = g \cdot f(x)$ for all $g \in G$ and $x \in X$. If the action of G on Y is trivial (i.e., $g \cdot y = y$ for all $g \in G$ and $y \in Y$), then f is called *invariant*, meaning $f(g \cdot x) = f(x)$ for all $g \in G$ and $x \in X$.

In the context of static graph representation learning, permutation equivariance ensures that the model’s output does not depend on the arbitrary ordering of nodes. Additionally, when modelling time-dependent data, such as in temporal graphs, it is also natural to require equivariance under reparametrisations of the time axis. We now formalise both types of equivariance:

1. **Permutation equivariance:** Given a permutation $p : [n] \rightarrow [n]$ with corresponding matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$, a function $f : \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times d}$ is *permutation equivariant* if $f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^T) = \mathbf{P}f(\mathbf{X}, \mathbf{A})$ for all permutations p and all graphs with node features \mathbf{X} and adjacency matrix \mathbf{A} .

2. **Time-warp equivariance:** Let $X : [0, T] \rightarrow Y$ be a continuous path. A *time-warp* is a diffeomorphism (a smooth bijection with smooth inverse) $\tau : [0, T] \rightarrow [0, T]$ satisfying $\tau(0) = 0$ and $\tau(T) = T$. A function $f : \chi([0, T], Y) \rightarrow \chi([0, T], Y)$ is *time-warp equivariant* if $f(X \circ \tau) = f(X) \circ \tau$ for all time-warps τ and paths X .

For a more formal treatment of group actions and equivariance, we refer the reader to Appendix C.

2.2 Neural differential equations for time-series data

Neural Differential Equations (NDEs) ([46, 51, 7, 34]) have emerged as a powerful tool at the intersection of dynamical systems and deep learning. In [8], the authors viewed the layers of a neural network as the time variable of an Ordinary Differential Equation (ODE), introducing the Neural Ordinary Differential Equation (NODE). In NODEs, the time dimension is merely an internal detail of the model, and the trajectories z are entirely determined by the initial condition. To extend the NODE framework to sequential data, [36] introduced *Neural Controlled Differential Equations* (NCDEs) via controlled differential equations (CDEs). For a continuous driving path $X : [0, T] \rightarrow \mathbb{R}^{d_x}$ NCDEs learn a latent path $z : [0, T] \rightarrow \mathbb{R}^{d_z}$ via the CDE

$$z(0) = \ell_\theta^1(X(0)), \quad z(t) = z(0) + \int_0^t f_\theta(z(s))dX(s) \quad (1)$$

which then returns either a scalar output $y \approx \ell_\theta^2(z(T))$ or an output path $y(t) \approx \ell_\theta^2(z(t))$. Here, $dX(s)$ denotes the Riemann-Stieltjes integral. In practice, we often only have discrete observations $((t_0, X_{t_0}), \dots, (t_N, X_{t_N}))$ with $t_j \in \mathbb{R}$ and $X_{t_j} \in \mathbb{R}^{d_x}$. To address this, we interpolate the observations into a continuous driving path $X : [t_0, t_N] \rightarrow \mathbb{R}^{d_x+1}$ such that $X(t_j) = (t_j, X_{t_j})$ for all j . If the path X is differentiable and has bounded derivative, the Riemann-Stieltjes integral, and hence NCDEs, can be rewritten as the following ordinary integral:

$$z(t) = z(0) + \int_0^t f_\theta(z(s)) \frac{dX(s)}{ds} ds \quad \text{for } t \in (0, T]. \quad (2)$$

The choice of interpolation scheme can impact the performance of an NCDE, with [45] providing a discussion on the theoretical properties and practical performance of a range of choices. Additionally, Log-NCDEs [63] extend NCDEs to non-differentiable paths X by leveraging the Log-ODE method to approximate solutions to Equation (1). Empirically, this improves training stability and model performance, especially for long time-series.

2.3 Permutation invariant and equivariant linear functions

Graph Neural Networks are typically constructed as permutation equivariant functions by propagating information locally on graphs following a message passing paradigm. An alternative approach, proposed by [44], represents d -dimensional graph data on k -tuples of nodes as a single matrix $\mathbf{Y} \in \mathbb{R}^{n^k \times d}$. For example, the case $k = 1$ corresponds to node signals and $k = 2$ to signals on edges.

In their work, the authors provide a full characterisation of linear maps $L : \mathbb{R}^{n^k \times d} \rightarrow \mathbb{R}^{n^k \times d'}$ that are equivariant with respect to permutations of the underlying node set. Remarkably, they show the dimension of the vector space formed by these equivariant maps, denoted by $\mathfrak{E}_{\Sigma_n}(k, d)^{d'}$, depends only on k , d , and d' , and is independent of n . In Section 3, we will be interested in the basis of $\mathfrak{E}_{\Sigma_n}(2, 1)^1$, the vector space of linear maps $L : \mathbb{R}^{n^2} \rightarrow \mathbb{R}^{n^2}$ between edge-valued data, which has a dimension of 15. For a list of terms spanning this basis, see Appendix A.

3 Permutation equivariant neural graph controlled differential equations

This section brings together Neural CDEs and temporal graph representation learning. We begin by introducing a recent approach proposed in [48], called *Graph Neural Controlled Differential Equations*. Through the lens of Geometric Deep Learning [5], we will identify a key theoretical limitation of this approach and propose an improved solution.

3.1 Graph Neural Controlled Differential Equations

The concept of neural controlled differential equations has been extended to graphs by incorporating dynamic adjacency matrices to drive the CDE dynamics [48].

Graph Neural CDEs. Let $\zeta_\theta : \mathbb{R}^{n \times n \times 2} \rightarrow \mathbb{R}^{n \times d_z}$ and $f_\theta : \mathbb{R}^{n \times d_z} \times \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{(n \times d_z) \times (n \times n \times 2)}$ be two graph neural networks. The *neural controlled differential equation for dynamic graphs* is defined as

$$\mathbf{Z}_t = \mathbf{Z}_{t_0} + \int_{t_0}^t f_\theta(\mathbf{Z}_s, \mathbf{A}_s) d\hat{\mathbf{A}}_s \quad \text{for } t \in (t_0, t_N] \quad (3)$$

where $\mathbf{Z}_{t_0} = \zeta_\theta(\hat{\mathbf{A}}_{t_0})$ is the initial condition, $\hat{\mathbf{A}} : [t_0, t_N] \rightarrow \mathbb{R}^{n \times n \times 2}$ such that $\hat{\mathbf{A}}_{t_k}^{i,j} = (t_k, \mathbf{A}_{t_k}^{i,j})$, and the product $f_\theta(\mathbf{Z}_s, \mathbf{A}_s) d\hat{\mathbf{A}}_s$ is a tensor contraction over $\mathbb{R}^{n \times n \times 2}$. The final prediction $\tilde{\mathbf{Y}}_t \in \mathbb{R}^{n \times d_y}$ is attained by row-wise application of another linear function $\ell_\theta : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_y}$, i.e. by slight abuse of notation $\tilde{\mathbf{Y}}_t = \ell_\theta(\mathbf{Z}_t)$.

Practical Considerations. Implementing Equation 3 directly is computationally intractable due to the large output dimension of f_θ and the complexity of the tensor contraction under the integral sign. Hence, [48] proposes the following simplification based on a message passing paradigm:

$$\mathbf{Z}_t = \mathbf{Z}_{t_0} + \int_{t_0}^t \mathbf{Z}_s^{(L)} ds \quad \text{for } t \in (t_0, t_N] \quad (4)$$

where $\mathbf{Z}_s^{(l)} = \sigma(\tilde{\mathbf{A}}_s \mathbf{Z}_s^{(l-1)} \mathbf{W}^{(l-1)})$ for $l \in \{1, \dots, L\}$, the adjacency matrix and its derivative are fused via $\tilde{\mathbf{A}}_s = \mathbf{W}^{(F)} \left[\frac{\mathbf{A}_s}{d\mathbf{A}_s} \right]$ with $\mathbf{W}^{(F)} \in \mathbb{R}^{n \times 2n}$ being a learnable fusion matrix, and $\mathbf{Z}_s^{(0)} = \mathbf{Z}_s$. Importantly, this simplification preserves the multiplicative interaction between the hidden state and control path which are critical for expressivity [13].

3.2 Inducing permutation equivariance

As motivated in Section 2, it would be natural to require equivariance with respect to permutation of the node and edge sets for any function on graph data. However, neither the original GN-CDE formulation in Equation 3, nor its approximation in Equation 4 satisfy this property. See Appendix E for a detailed proof.

In the following, we introduce a fully permutation equivariant variant of the GN-CDE framework. By leveraging the characterisation of linear permutation equivariant layers presented in Section 2.3, our goal is to approximate Equation 4 while maintaining equivariance. Since the equivariance breaks in the fusion step, our strategy is to project the fusion operation onto the subspace of linear equivariant functions. Specifically, if we interpret the fusion as the application of linear maps $L_1, L_2 : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ to the matrices \mathbf{A}_s and $\frac{d\mathbf{A}_s}{ds}$, respectively, we can project L_1 and L_2 onto $\mathfrak{E}_{\Sigma_n}(2, 1)^1$. According to the Projection Theorem (see Appendix B), this yields the best possible approximation. This motivates us to propose the following:

The model. Let $L_1, L_2 \in \mathfrak{E}_{\Sigma_n}(2, 1)^1$ be two learnable permutation equivariant linear maps, i.e. a weighted combination of the basis terms in Appendix A. Using these, we combine the adjacency and its time derivative as

$$\bar{\mathbf{A}}_s = L_1(\mathbf{A}_s) + L_2\left(\frac{d\mathbf{A}_s}{ds}\right). \quad (5)$$

Now, let σ be a non-linear activation function and denote by $\mathbf{Z}_s^{(L)}$ the latent representation obtained by an iterative convolution operation of the form $\mathbf{Z}_s^{(l)} = \sigma(\bar{\mathbf{A}}_s \mathbf{Z}_s^{(l-1)} \mathbf{W}^{(l-1)})$ for $l \in \{1, \dots, L\}$ where the $\mathbf{W}^{(l)}$ are learnable matrices. The *Permutation Equivariant Neural Graph Controlled Differential Equation (PENG-CDE)* then takes the form

$$\mathbf{Z}_t = \mathbf{Z}_{t_0} + \int_{t_0}^t \sigma(\bar{\mathbf{A}}_s \mathbf{Z}_s^{(L)} \mathbf{W}^{(L)}) ds \quad \text{for } t \in (t_0, t_N]. \quad (6)$$

with initial condition $\mathbf{Z}_{t_0} = \zeta_\theta(\hat{\mathbf{A}}_{t_0})$.

Theoretical Properties. In the purely linear case (i.e., without applying the non-linearities σ in between the layers), one can show that the notion of optimality above, as motivated by the Projection Theorem, is satisfied, as formalised in the following theorem:

Theorem 3.1. *In the absence of non-linearities, the PENG-CDE model in Equation 6 is the projection of the model in Equation 4 onto the space of equivariant linear functions.*

Proof sketch. Projecting an the flow of an ODE onto a function space is equivalent to projecting its associated vector field. Moreover, note that the vector field in Equation 4 can be decomposed as the composition of a standard convolutional graph neural network, which is inherently equivariant, and the adjacency fusion, which is not. Thus, projecting Equation 4 is equivalent to projecting only the fusion operator. For a detailed proof, see Appendices C and D. \square

Although the projection-based notion of optimality holds strictly only in the linear case, by construction, our model satisfies both equivariance constraints outlined in Section 2.1.

Proposition 3.2. *The PENG-CDE model in Equation 6 is both permutation equivariant in the spatial domain and time-warp equivariant in the temporal domain.*

Proof. See Appendix E. \square

3.3 Including dynamic node features

So far, our framework only incorporates dynamic adjacency matrices into the latent dynamics, leaving dynamic node features unaddressed. To remedy this, we draw on the approach of [11]. Concretely, suppose node feature snapshots $\{X_{t_k}\}_{k=0}^N$ are available, where each $X_{t_k} \in \mathbb{R}^{d_x}$. We interpolate these snapshots to obtain a continuous, differentiable path $\mathbf{X} : [t_0, t_N] \rightarrow \mathbb{R}^{d_x+1}$ as in Section 2.2. Next, we choose $\mathbf{W}^{(L)} \in \mathbb{R}^{d_z \times (d_z \times (d_x+1))}$ to set the output dimension of $\sigma(\bar{\mathbf{A}}_s \mathbf{Z}_s^{(L)} \mathbf{W}^{(L)})$ to $\mathbb{R}^{n \times (d_z \times (d_x+1))}$, enabling a row-wise (i.e., node-wise) Hadamard multiplication, denoted by \odot . The resulting system is

$$\mathbf{Z}_t = \mathbf{Z}_{t_0} + \int_{t_0}^t \left\{ \sigma(\bar{\mathbf{A}}_s \mathbf{Z}_s^{(L)} \mathbf{W}^{(L)}) \odot \frac{d\mathbf{X}_s}{ds} \right\} ds \quad \text{for } t \in (t_0, t_N]. \quad (7)$$

Building on Proposition 3.2, this extension preserves both permutation equivariance in the spatial domain and time-warp equivariance in the temporal domain. Intuitively, this is because a permutation of the node indices reorders both $\sigma(\bar{\mathbf{A}}_s \mathbf{Z}_s^{(L)} \mathbf{W}^{(L)})$ and $\frac{d\mathbf{X}_s}{ds}$ in the same manner, ensuring that the overall model remains equivariant. A detailed proof is provided in Appendix E. \square

4 Numerical experiments

In this section, we evaluate the PENG-CDE model on a range of synthetic and real-world tasks. First, we replicate the experiments from [48] and compare them with the non-permutation-equivariant version. Next, we evaluate the model on well-established real-world dynamic graph benchmarks against other common approaches. Finally, we conduct an ablation study to examine the weighting of different basis terms of $\mathfrak{E}_{\Sigma_n}(2, 1)^1$ in Appendix H. For supplementary information regarding implementation details, see Appendix F.

4.1 Synthetic experiments: heat diffusion and gene regulation

Task. We randomly sample initial graphs from four distinct graph distributions (grid, small-world, power-law, and community), each comprising 400 nodes. We then take 120 irregularly sampled time-stamps spanning $T = 0$ to $T = 5$. At 12 time steps, sampled uniformly at random from these snapshots, we randomly add or remove edges following a Bernoulli trial. The final 20 snapshots are allocated for extrapolation validation, while from the remaining 100, a random subset of 20 is used for interpolation validation and the remaining 80 for training. A batch of four such time series are generated for each of training, validation, and testing. We then simulate node features according to the heat diffusion dynamics governed by Newton’s law of cooling and the gene regulatory dynamics governed by the Michaelis–Menten equation.

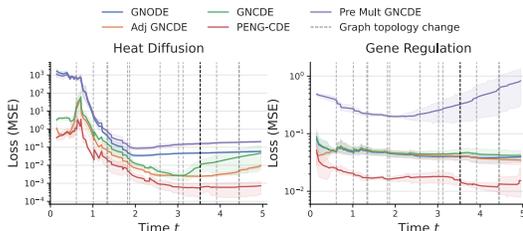


Figure 1: Test losses, plotted against simulation time, for the Graph Neural ODE, three GN-CDE variants, and our proposed Permutation-Equivariant GN-CDE model on the heat diffusion (left) and gene regulation (right) tasks. Dashed vertical lines mark changes in graph topology, while the bold black line indicates the final time point in the training set. Results are reported as means (solid) and ranges (shaded) over a test set with a batch size of four.

adjacency matrix and its derivative (i.e. $\tilde{\mathbf{A}}_s = \mathbf{A}_s + \frac{d\mathbf{A}_s}{ds}$), which implicitly achieves permutation equivariance - a property not present in the original formulation. To ensure a fair comparison with a non-equivariant GN-CDE variant, we also consider the Premultiplication Fusion GN-CDE (Pre Mult GN-CDE) model. In this variant, the fusion step is performed via premultiplication by learnable matrices \mathbf{W}_1 and \mathbf{W}_2 , such that $\tilde{\mathbf{A}}_s = \mathbf{W}_1\mathbf{A}_s + \mathbf{W}_2\frac{d\mathbf{A}_s}{ds}$. Lastly, we include our proposed Permutation Equivariant Neural Graph CDE (PENG-CDE). Experimental results for the heat diffusion and gene regulation tasks are summarised in Table 1. Moreover, we visualise the test loss over simulation time for models initialised from the community graph distribution in Figure 1.

Finding I: Equivariance improves performance. Incorporating permutation equivariance (found in models Adjacency GN-CDE, Original GN-CDE, and PENG-CDE) leads to performance improvements of an order of magnitude over the non-equivariant Pre Mult GN-CDE across all considered tasks and graph distributions.

Finding II: Enhanced expressivity via 15 basis terms boosts performance. Integrating all 15 basis terms of linear equivariant maps into the fusion of our model significantly enhances its expressivity. Compared to the original implementation, the PENG-CDE achieves relative MSE

Baselines. We compare the performance of the following models. As a simple baseline, we first consider a model with a constant vector field (Const), that is $\mathbf{Z}_s^{(L)} = \mathbf{b}$ for all $s \in [t_0, t_N]$ in Equation 4. Next, we consider the Graph Neural ODE (GNODE) model [47], in which the adjacency matrix within the vector field is discretised by flooring the time index. The GN-ODE is governed by the equation:

$$\mathbf{Z}_t = \mathbf{Z}_{t_0} + \int_{t_0}^t f_{\theta}(\mathbf{Z}_s, \mathbf{A}_{\lfloor s \rfloor}) ds \quad (8)$$

for all $t \in (t_0, t_N]$ where $\mathbf{A}_{\lfloor s \rfloor}$ is the adjacency matrix corresponding to the graph \mathcal{G}_{t_k} with $t_k = \lfloor s \rfloor$. Additionally, we consider several variants of the GN-CDE. Firstly, we include a simplified fusion model, named Adjacency GN-CDE, which employs interpolated adjacency matrices, i.e. $\tilde{\mathbf{A}}_s = \mathbf{A}_s$. We note that for the experiments on these two datasets, the GN-CDE [48] model implemented the fusion by a simple element-wise summation between the

improvements ranging from 30.44% to 73.84% on the heat diffusion task and from 39.71% to 67.06% on gene regulation tasks. One plausible explanation is that summing across rows (basis term 3 in Appendix A) corresponds to computing node degrees in an unweighted graph, which facilitates degree normalisation as required in Equation 24.

Finding III: Additional terms enhance extrapolation behaviour. As shown in Figure 1, for the heat diffusion task, our PENG-CDE is the only model capable of maintaining constant losses during the extrapolation phase (i.e. over the final 20 snapshots). This suggests that the added equivariant terms contribute to more robust extrapolation, likely due to mechanisms similar to those discussed above.

Table 1: Comparison of GN-CDE variants and baselines on the heat diffusion (top) and gene regulation (bottom) tasks. Mean MSEs with 95% confidence intervals are reported, with the best mean highlighted in **bold**, and all results within the corresponding confidence interval are underlined. The final row in each table reports the relative improvement of PENG-CDE over the original GN-CDE formulation.

Heat Diffusion Task (MSE ↓)				
Model	Community	Grid	Power Law	Small World
Constant	1.936 ± 0.550	11.155 ± 0.669	3.147 ± 0.850	6.286 ± 1.056
Graph Neural ODE [47]	0.237 ± 0.322	1.001 ± 0.751	<u>0.270 ± 0.310</u>	<u>0.311 ± 0.268</u>
Adjacency GN-CDE	0.208 ± 0.240	0.691 ± 0.887	0.258 ± 0.288	<u>0.248 ± 0.240</u>
Pre Mult GN-CDE	1.968 ± 0.548	12.262 ± 1.437	7.440 ± 4.458	6.829 ± 0.644
Original GN-CDE [48]	0.366 ± 0.400	1.324 ± 0.630	0.417 ± 0.314	0.552 ± 0.470
PENG-CDE (ours)	0.096 ± 0.051	0.481 ± 0.195	<u>0.290 ± 0.265</u>	0.247 ± 0.215
Relative Improvement	73.84%	63.70%	30.44%	55.20%

Gene Regulation Task (MSE ↓)				
Model	Community	Grid	Power Law	Small World
Constant	36.307 ± 2.609	1.390 ± 0.168	7.099 ± 0.382	0.772 ± 0.126
Graph Neural ODE [47]	8.548 ± 3.212	0.167 ± 0.191	0.372 ± 0.398	0.294 ± 1.308
Adjacency GN-CDE	8.909 ± 6.311	1.476 ± 2.504	0.498 ± 0.126	0.195 ± 0.046
Pre Mult GN-CDE	153.084 ± 149.609	2.553 ± 0.251	6.978 ± 3.045	1.591 ± 0.146
Original GN-CDE [48]	10.717 ± 7.079	0.457 ± 0.167	0.822 ± 0.299	0.323 ± 0.154
PENG-CDE (ours)	4.566 ± 2.780	<u>0.247 ± 0.090</u>	<u>0.526 ± 0.220</u>	0.186 ± 0.484
Relative Improvement	67.06%	45.90%	39.71%	54.14%

Oversampling and irregularity.

In Figure 2, we compare the performance of the PENG-CDE with DCRNN [40], STIDGCN [41] and ASTGCN [26] on data generated from the graph SIR disease spread model [33]. Its parameters are chosen to produce trajectories both with and without an outbreak, and the task is to classify each trajectory into one of the two categories. To study the effect of oversampling, we simulate each trajectory up to a fixed terminal time $T = 1$ and train and evaluate the models on datasets with an increasing number of observation points (left and middle panels). We also investigate how performance varies with the regularity of the sampling grid: by drawing observation times from a Gamma distribution with controlled shape parameter k , we obtain series with different degrees of irregularity and plot test accuracy as irregularity increases (right panel). For implementation details, see Appendix F.3.

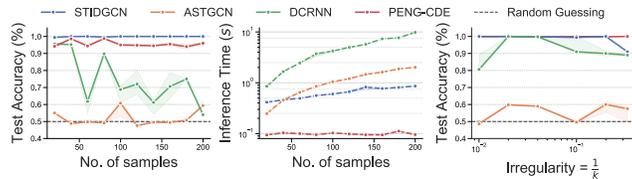


Figure 2: Classification accuracy on the SIR model for STIDGCN, ASTGCN, DCRNN, and our PENG-CDE as a function of the number of observed timesteps (left) and sampling irregularity (right). Inference time with increasing numbers of observations is shown in the middle panel.

Table 2: Results of experiments for the `england-covid` and `twitter-tennis` tasks from the Pytorch Geometric Temporal datasets [54]. Mean values and 95% confidence intervals are reported, with the best mean highlighted in **bold** and all results within the confidence interval around the best mean are underlined.

Model	england-covid		twitter-tennis	
	MSE ↓		MSE ↓	
	Validation	Test	Validation	Test
DCRNN [40]	<u>0.898 ± 0.126</u>	<u>1.021 ± 0.203</u>	<u>0.472 ± 0.136</u>	<u>0.455 ± 0.087</u>
ASTGCN [26]	1.235 ± 0.139	1.283 ± 0.257	0.545 ± 0.128	0.530 ± 0.078
STIDGCN [41]	0.916 ± 0.111	0.933 ± 0.100	0.524 ± 0.140	0.495 ± 0.059
GNODE [47]	0.802 ± 0.184	<u>0.945 ± 0.251</u>	<u>0.419 ± 0.068</u>	<u>0.525 ± 0.043</u>
STG-NCDE [11]	1.484 ± 0.515	<u>1.778 ± 1.084</u>	<u>0.372 ± 0.108</u>	<u>0.453 ± 0.035</u>
GN-CDE [48]	<u>0.892 ± 0.143</u>	<u>0.962 ± 0.278</u>	0.369 ± 0.070	<u>0.443 ± 0.053</u>
PENG-CDE (ours)	<u>0.836 ± 0.122</u>	0.913 ± 0.200	<u>0.391 ± 0.069</u>	0.440 ± 0.053

Finding IV: PENG-CDE is robust to oversampling and irregular sampling. CDE-based models like the PENG-CDE decouple the complexity of their forward passes from the number of input observations: the number of vector field evaluations depends on the ODE solver rather than the sampling rate. This makes them inherently more robust to oversampling. As shown in Figure 2, increasing the number of observations and the level of irregularity degrades the performance of the recurrent model, while PENG-CDE retains stable accuracy. Although STIDGCN also maintains performance under oversampling, its computational cost increases with the number of observations, in contrast to the constant runtime of our approach.

4.2 Real-world tasks

We evaluate our framework on two types of real-world datasets: snapshot-based and event-based. First, we consider two node regression tasks from PyTorch Geometric Temporal [54] and then two node affinity prediction tasks from the Temporal Graph Benchmark (TGB) [29, 24]. Implementation details are provided in Appendix F.4.

Snapshot-based datasets. In this setting, we observe a sequence of dense graphs, one at each timestep. To demonstrate that our modifications enable our model to better capture dynamic graph topologies, we compare it against several differential equation-based models: the Graph Neural ODE [47], the Spatio-Temporal Graph Neural Controlled Differential Equation (STG-NCDE) [11], and the original GN-CDE [48]. Additionally, we include comparisons with recurrent (DCRNN [40]), attentional (ASTGCN [26]), and interactive (STIDGCN [41]) models. We benchmark performance using the `england-covid` and `twitter-tennis` datasets from PyTorch Geometric Temporal [54] - the only datasets in the library exhibiting dynamic graph topologies. Table 2 reports the mean MSE with 95% confidence intervals on both the validation and test sets over 10 random seeds.

Finding V: PENG-CDE generalises well on dynamic snapshot datasets. While not achieving the lowest validation error, PENG-CDE attains the lowest test error on both tasks, indicating better generalisation and reduced overfitting compared to other models.

Event-based datasets. In many applications, a temporal graph is represented as a sequence of events $\mathcal{G} = \{(e_i, t_i, x_i)\}_{i=1}^n$, where each event e_i (e.g., an interaction between nodes u and v) occurs at time t_i and is associated with data $x_i \in \mathbb{R}^d$. Since GN-CDEs are not inherently designed for processing individual events, we aggregate events into snapshots. In more detail, given a time window of length Δt , we partition the overall time span $[t_0, t_n]$ into $m = \lfloor \frac{t_n - t_0}{\Delta t} \rfloor$ non-overlapping intervals. For each interval indexed by $k \in \{0, 1, \dots, m-1\}$, we define the snapshot graph \mathcal{G}_k to include all events (or edges) that occur within the time window $[t_0 + k\Delta t, t_0 + (k+1)\Delta t]$. This then gives us a sequence of snapshots $\{\mathcal{G}_0, \dots, \mathcal{G}_{m-1}\}$, which we can integrate into the setting above. We evaluate our framework in this event-based setting on the trade and genre node affinity prediction tasks from the Temporal Graph Benchmark (TGB) [29, 24]. These two tasks were selected from the four available in the benchmark, as the dense adjacency matrices required by Neural ODE-based models make the remaining tasks computationally infeasible within our experimental constraints.

Table 3: Experimental results for the `tgbn-trade` and `tgbn-genre` node affinity prediction tasks from the Temporal Graph Benchmark datasets [29, 24]. Results marked with † are from [29], those with ‡ are from [67], and those with * are from [58]. Mean values and standard deviations are reported with results within one standard deviation of the best mean highlighted in **bold** (for deterministic and learned models separately).

Model	trade NDCG@10 †	genre NDCG@10 †
Persistent Forecast (L) †	0.855	0.357
Moving Avg (L) †	0.823	0.509
Moving Avg (M)	0.777	0.472
JODIE ‡ [39]	0.374±0.09	0.350±0.04
TGAT ‡ [15]	0.375±0.07	0.352±0.03
CAWN ‡ [65]	0.374±0.09	–
TCL ‡ [64]	0.375±0.09	0.354±0.02
GraphMixer ‡ [14]	0.375±0.11	0.352±0.03
DyGFormer ‡ [67]	0.388±0.64	0.365±0.20
DyRep † [59]	0.374±0.001	0.351±0.001
TGN † [52]	0.374±0.001	0.367±0.058
TGNv2* [58]	0.735±0.006	0.469±0.002
STG-NCDE [10]	0.618±0.024	0.438±0.038
GN-CDE [48]	0.713±0.026	0.460±0.016
PENG-CDE	0.716±0.029	0.523±0.017
+ Source/Target Id	0.734±0.024	–

Table 3.

Finding VI: PENG-CDE achieves state-of-the-art performance on TGB. Without source-target identification, our model is outperformed only by TGNv2 on the `tgbn-trade` task; with source-target identification, it performs on par. Moreover, on the `tgbn-genre` task, our model significantly outperforms all other models and is the only machine-learned approach to surpass the performance of all heuristic baselines.

5 Conclusion

In this work, we introduced Permutation Equivariant Neural Graph CDEs - a geometrically grounded approach to temporal graph representation learning that balances parameter efficiency with formal expressivity. In the linear setting, we showed that our model can be derived as a projection of Graph Neural CDEs onto the space of permutation equivariant functions. Through synthetic experiments, we demonstrated that both the imposed equivariance and enhanced expressivity improve performance in both interpolation and extrapolation tasks. Additionally, our model retains the strengths of Neural CDEs in handling oversampling and irregular time intervals. On the TGB-`genre` dataset, our model achieves a new state-of-the-art - becoming the first learned method to surpass a moving average baseline.

5.1 Limitations and future work

A current limitation of our framework and implementation is their reliance on dense adjacency matrices, which restricts scalability to large, sparse graphs. Future work includes addressing this scalability bottleneck and investigating the impact of more tailored hyperparameter configurations, such as the choice of interpolation schemes, ODE solvers, and vector field architectures. In addition, advanced solvers like Log-NCDE [63], originally developed for Euclidean paths, could be adapted to the graph setting.

Baselines. For the event-based experiments, we first compare against three simple heuristics based on ground-truth labels and messages: ‘Persistent Forecast (L)’ and ‘Moving Average (L)’, which operate on labels, and ‘Moving Average (M)’, which is applied to messages. We then benchmark several learned models, including JODIE [39], TGAT [15], CAWN [65], TCL [64], GraphMixer [14], DyGFormer [67], DyRep [59], TGN [52], and TGNv2 [58]. Additionally, we include STG-NCDE [10], GN-CDE [48], and our proposed PENG-CDE. It is a known issue that the heuristics above often outperform learned models on node-affinity prediction tasks [24]; see [58] for a detailed discussion. This motivated TGNv2 to introduce a mechanism that learns node embeddings to distinguish source and target nodes for each interaction, a feature tailored specifically to the node affinity prediction task. To ensure a fair comparison, we incorporate an analogous mechanism into a variant of our model. Further details are provided in Appendix F.5. The experimental results are summarised in

While we chose dynamic adjacency matrices as the control signal for the CDE, this is just one possible vectorisation of graph structure. Prior work [25] suggests that adjacency matrices are often suboptimal for representing graph properties. Identifying more expressive or task-aligned graph representations could further improve performance and remains an exciting research direction.

Finally, while Neural CDEs are known to be universal approximators on Euclidean domains [36], and recent work has also provided generalisation bounds [3], no such theoretical guarantees currently exist for graph-based Neural CDEs. Establishing these would be a valuable theoretical contribution to the growing field of temporal graph representation learning.

References

- [1] J.P. Aubin. *Applied Functional Analysis*. A Wiley-Interscience publication. Wiley, 1979. ISBN 9780471021490. URL <https://books.google.co.uk/books?id=CQzE1g0HzP4C>.
- [2] Richard Bergna, Sergio Calvo-Ordóñez, Felix Opolka, Pietro Liò, and Jose Miguel Hernandez-Lobato. Uncertainty modeling in graph neural networks via stochastic differential equations, 2025.
- [3] Linus Bleistein and Agathe Guilloux. On the Generalization Capacities of Neural Controlled Differential Equations. In *Workshop on New Frontiers in Learning, Control, and Dynamical Systems at the International Conference on Machine Learning (ICML 2023)*, Honolulu (Hawaii), United States, July 2023. URL <https://hal.science/hal-04876009>.
- [4] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- [5] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges, 2021.
- [6] Ben Chamberlain, James Rowbottom, Maria I. Gorinova, Michael Bronstein, Stefan Webb, and Emanuele Rossi. GRAND: Graph Neural Diffusion. In *Proceedings of the 38th International Conference on Machine Learning*, pages 1407–1418. PMLR, July 2021. URL <https://proceedings.mlr.press/v139/chamberlain21a.html>. ISSN: 2640-3498.
- [7] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, page 6572–6583, Red Hook, NY, USA, 2018. Curran Associates Inc.
- [8] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/69386f6bb1dfed68692a24c8686939b9-Paper.pdf.
- [9] Kyunghyun Cho, B van Merriënboer, Caglar Gulcehre, F Bougares, H Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014.
- [10] Jeongwhan Choi and Noseong Park. Graph neural rough differential equations for traffic forecasting. *ACM Trans. Intell. Syst. Technol.*, 14(4), July 2023. ISSN 2157-6904. doi: 10.1145/3604808. URL <https://doi.org/10.1145/3604808>.
- [11] Jeongwhan Choi, Hwangyong Choi, Jeehyun Hwang, and Noseong Park. Graph neural controlled differential equations for traffic forecasting. In *AAAI*, 2022.
- [12] Jeongwhan Choi, Seoyoung Hong, Noseong Park, and Sung-Bae Cho. Gread: Graph neural reaction-diffusion networks. In *ICML*, 2023.
- [13] Nicola Muca Cirone, Antonio Orvieto, Benjamin Walker, Cristopher Salvi, and Terry Lyons. Theoretical foundations of deep selective state-space models. *Advances in Neural Information Processing Systems*, 2024.
- [14] Weilin Cong, Si Zhang, Jian Kang, Baichuan Yuan, Hao Wu, Xin Zhou, Hanghang Tong, and Mehrdad Mahdavi. Do we really need complicated model architectures for temporal networks?, 2023. URL <https://arxiv.org/abs/2302.11636>.
- [15] da Xu, chuanwei ruan, evren korpeoglu, sushant kumar, and kannan achan. Inductive representation learning on temporal graphs. In *International Conference on Learning Representations (ICLR)*, 2020.

- [16] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS'16, page 3844–3852, Red Hook, NY, USA, 2016. Curran Associates Inc. ISBN 9781510838819.
- [17] Moshe Eliasof, Eldad Haber, and Eran Treister. Graph Neural Reaction Diffusion Models. 46(4):C399–C420. ISSN 1064-8275, 1095-7197. doi: 10.1137/23M1576700. URL <https://epubs.siam.org/doi/10.1137/23M1576700>.
- [18] Moshe Eliasof, Eldad Haber, and Eran Treister. Pde-gcn: novel architectures for graph neural networks motivated by partial differential equations. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- [19] Moshe Eliasof, Eldad Haber, and Eran Treister. Feature transportation improves graph neural networks. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI'24/IAAI'24/EAAI'24. AAAI Press, 2024. ISBN 978-1-57735-887-9. doi: 10.1609/aaai.v38i11.29073. URL <https://doi.org/10.1609/aaai.v38i11.29073>.
- [20] Moshe Eliasof, Eldad Haber, Eran Treister, and Carola-Bibiane Sch"onlieb. Data-driven higher order differential equations inspired graph neural networks. In *ICLR 2024 Workshop on AI4DifferentialEquations In Science*, 2024. URL <https://openreview.net/forum?id=rJReXWFBYt>.
- [21] Moshe Eliasof, Eldad Haber, Eran Treister, and Carola-Bibiane B Sch"onlieb. On the temporal domain of differential equation inspired graph neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 1792–1800. PMLR, 2024.
- [22] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. Graph neural networks for social recommendation. In *The World Wide Web Conference*, WWW '19, page 417–426, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366748. doi: 10.1145/3308558.3313488. URL <https://doi.org/10.1145/3308558.3313488>.
- [23] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [24] Julia Gastinger, Shenyang Huang, Mikhail Galkin, Erfan Loghmani, Ali Parviz, Farimah Pour-safaei, Jacob Danovitch, Emanuele Rossi, Ioannis Koutis, Heiner Stuckenschmidt, Reihaneh Rabbany, and Guillaume Rabusseau. Tgb 2.0: A benchmark for learning on temporal knowledge graphs and heterogeneous graphs. *Advances in Neural Information Processing Systems*, 2024.
- [25] Martin Grohe. word2vec, node2vec, graph2vec, x2vec: Towards a theory of vector embeddings of structured data. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS'20, page 1–16, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450371087. doi: 10.1145/3375395.3387641. URL <https://doi.org/10.1145/3375395.3387641>.
- [26] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention Based Spatial-Temporal Graph Convolutional Networks for Traffic Flow Forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):922–929, July 2019. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v33i01.3301922.
- [27] Shengnan Guo, Youfang Lin, Huaiyu Wan, Xiucheng Li, and Gao Cong. Learning Dynamics and Heterogeneity of Spatial-Temporal Graph Data for Traffic Forecasting. *IEEE Transactions on Knowledge and Data Engineering*, 34(11):5415–5428, November 2022. ISSN 1041-4347, 1558-2191, 2326-3865. doi: 10.1109/TKDE.2021.3056502.
- [28] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith,

- Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- [29] Shenyang Huang, Farimah Poursafaei, Jacob Danovitch, Matthias Fey, Weihua Hu, Emanuele Rossi, Jure Leskovec, Michael Bronstein, Guillaume Rabusseau, and Reihaneh Rabbany. Temporal graph benchmark for machine learning on temporal graphs. *Advances in Neural Information Processing Systems*, 2023.
- [30] Weiwei Jiang and Jiayun Luo. Graph neural network for traffic forecasting: A survey. 207: 117921. ISSN 09574174. doi: 10.1016/j.eswa.2022.117921. URL <https://linkinghub.elsevier.com/retrieve/pii/S0957417422011654>.
- [31] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2. URL <https://doi.org/10.1038/s41586-021-03819-2>.
- [32] W.O. Kermack and A.G. McKendrick. A contribution to the mathematical theory of epidemics. 115(772):700–721. ISSN 0950-1207, 2053-9150. doi: 10.1098/rspa.1927.0118. URL <https://royalsocietypublishing.org/doi/10.1098/rspa.1927.0118>.
- [33] C. Kerr. The Mathematical Theory of Infectious Diseases and its Applications. By Norman T. J. Bailey, M.A., D.Sc.; second edition, 1975. London: Charles Griffin & Co. Ltd. 9*x6*, pp. 430, with diagrams. Price: £14.00. 1(18):674. ISSN 0025-729X, 1326-5377. doi: 10.5694/j.1326-5377.1976.tb140951.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.5694/j.1326-5377.1976.tb140951.x>.
- [34] Patrick Kidger. *On Neural Differential Equations*. PhD thesis, University of Oxford, 2021.
- [35] Patrick Kidger and Cristian Garcia. Equinox: neural networks in JAX via callable PyTrees and filtered transformations. *Differentiable Programming workshop at Neural Information Processing Systems 2021*, 2021.
- [36] Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural Controlled Differential Equations for Irregular Time Series. *Advances in Neural Information Processing Systems*, 2020.
- [37] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations (ICLR)*, 2015.
- [38] Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *ArXiv*, abs/1609.02907, 2016. URL <https://api.semanticscholar.org/CorpusID:3144218>.
- [39] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2019.
- [40] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations (ICLR '18)*, 2018.
- [41] Aoyu Liu and Yaying Zhang. Spatial–Temporal Dynamic Graph Convolutional Network With Interactive Learning for Traffic Forecasting. *IEEE Transactions on Intelligent Transportation Systems*, 25(7):7645–7660, July 2024. ISSN 1524-9050, 1558-0016. doi: 10.1109/TITS.2024.3362145.

- [42] Shixuan Liu, Jijun Zhu, Wei Lei, and Ping Zhang. Spatial-Temporal Attention Graph WaveNet for Traffic Forecasting. In *2023 5th International Conference on Data-driven Optimization of Complex Systems (DOCS)*, pages 1–8. IEEE. ISBN 979-8-3503-9352-1. doi: 10.1109/DOCS60977.2023.10294485. URL <https://ieeexplore.ieee.org/document/10294485/>.
- [43] S. G. Lobanov and O. G. Smolyanov. Ordinary differential equations in locally convex spaces. *Uspekhi Mat. Nauk*, 39:93–168, 1994.
- [44] Haggai Maron, Heli Ben-Hamu, Nadav Shamir, and Yaron Lipman. Invariant and equivariant graph networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Syx72jC9tm>.
- [45] James Morrill, Patrick Kidger, Lingyi Yang, and Terry Lyons. On the choice of interpolation scheme for neural CDEs. *Transactions of Machine Learning Research*, 2022.
- [46] Pearlmutter. Learning state space trajectories in recurrent neural networks. In *International Joint Conference on Neural Networks*, pages 365–372 vol.2. IEEE. doi: 10.1109/IJCNN.1989.118724. URL <http://ieeexplore.ieee.org/document/118724/>.
- [47] Michael Poli, Stefano Massaroli, Junyoung Park, Atsushi Yamashita, Hajime Asama, and Jinkyoo Park. Graph neural ordinary differential equations. *CoRR*, abs/1911.07532, 2019. URL <http://arxiv.org/abs/1911.07532>.
- [48] Tiexin Qin, Benjamin Walker, Terry Lyons, Hong Yan, and Haoliang Li. Learning dynamic graph embeddings with neural controlled differential equations, 2023.
- [49] Jason Rader, Terry Lyons, and Patrick Kidger. Lineax: unified linear solves and linear least-squares in jax and equinox. *AI for science workshop at Neural Information Processing Systems 2023*, *arXiv:2311.17283*, 2023.
- [50] J. Rapin and J.-R. King. Exca - Execution and caching. <https://github.com/facebookresearch/exca>, 2024.
- [51] R. Rico-Martínez, K. Krischer, I.G. Kevrekidis, M.C. Kube, and J.L. Hudson. DISCRETE- vs. CONTINUOUS-TIME NONLINEAR SIGNAL PROCESSING OF Cu ELECTRODISSOLUTION DATA. 118(1):25–48. ISSN 0098-6445, 1563-5201. doi: 10.1080/00986449208936084. URL <https://www.tandfonline.com/doi/full/10.1080/00986449208936084>.
- [52] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. Temporal graph networks for deep learning on dynamic graphs. In *ICML 2020 Workshop on Graph Representation Learning*, 2020.
- [53] James Rowbottom, Benjamin P Chamberlain, Thomas Markovich, and Michael M Bronstein. Understanding convolution on graphs via energies. 2022.
- [54] Benedek Rozemberczki, Paul Scherer, Yixuan He, George Panagopoulos, Alexander Riedel, Maria Astefanoaei, Oliver Kiss, Ferenc Beres, Guzman Lopez, Nicolas Collignon, and Rik Sarkar. PyTorch Geometric Temporal: Spatiotemporal Signal Processing with Neural Machine Learning Models. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management*, page 4564–4573, 2021.
- [55] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009. doi: 10.1109/TNN.2008.2005605.
- [56] Alessandro Sperduti. Encoding labeled graphs by labeling raam. In J. Cowan, G. Tesauero, and J. Alspector, editors, *Advances in Neural Information Processing Systems*, volume 6. Morgan-Kaufmann, 1993. URL https://proceedings.neurips.cc/paper_files/paper/1993/file/fc49306d97602c8ed1be1dfbf0835ead-Paper.pdf.
- [57] Fang Sun, Zijie Huang, Haixin Wang, Huacong Tang, Xiao Luo, Wei Wang, and Yizhou Sun. Graph fourier neural odes: Modeling spatial-temporal multi-scales in molecular dynamics, 2025. URL <https://arxiv.org/abs/2411.01600>.

- [58] Benedict Aaron Tjandra, Federico Barbero, and Michael Bronstein. Enhancing the expressivity of temporal graph networks through source-target identification, 2024. URL <https://arxiv.org/abs/2411.03596>.
- [59] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. Representation learning over dynamic graphs, 2018. URL <https://arxiv.org/abs/1803.04051>.
- [60] Ch Tsitouras. Runge–kutta pairs of order 5 (4) satisfying only the first column simplifying assumption. *Computers & Mathematics with Applications*, 62(2):770–775, 2011.
- [61] Guido Van Rossum and Fred L Drake Jr. Python reference manual. *Centrum voor Wiskunde en Informatica Amsterdam*, 1995.
- [62] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. *6th International Conference on Learning Representations*, 2017.
- [63] Benjamin Walker, Andrew D. McLeod, Tiexin Qin, Yichuan Cheng, Haoliang Li, and Terry Lyons. Log neural controlled differential equations: The lie brackets make a difference. *International Conference on Machine Learning*, 2024.
- [64] Lu Wang, Xiaofu Chang, Shuang Li, Yunfei Chu, Hui Li, Wei Zhang, Xiaofeng He, Le Song, Jingren Zhou, and Hongxia Yang. Tcl: Transformer-based dynamic graph modelling via contrastive learning, 2021. URL <https://arxiv.org/abs/2105.07944>.
- [65] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. Inductive representation learning in temporal networks via causal anonymous walks, 2022. URL <https://arxiv.org/abs/2101.05974>.
- [66] Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, page 3634–3640. AAAI Press, 2018. ISBN 9780999241127.
- [67] Le Yu, Leilei Sun, Bowen Du, and Weifeng Lv. Towards better dynamic graph learning: New architecture and unified library, 2023. URL <https://arxiv.org/abs/2303.13047>.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: In the abstract and introduction, we make the following claims that reflect the paper's contributions and scope:

- a) Induce permutation equivariance into Graph Neural CDEs by projecting them onto permutation equivariant function spaces; formally prove this statement in the linear case
- b) Demonstrate theoretically that our model is equivariant with respect to permutations and time-warps
- c) Empirically validate our model on dynamic node-level tasks using synthetic experiments as well as standard benchmark datasets

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss limitations and future work in Section 5.1.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best

judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: All theoretical results in the main text (Theorem 3.1, Proposition 3.2) are proved in Appendix E. All new results in Appendices C and D are proved after their statement. All assumptions and auxiliary results are clearly stated.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide a full and detailed description of all experiments in the main text in Section 4, with additional details in Appendix F. An anonymised version of the source code used to run the experiments is provided at https://anonymous.4open.science/r/perm_equiv_gn_cdes-BBF8/README.md.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The TGB and Pytorch Geometric Temporal datasets used in Section 4 are publicly available. The data generation procedures for the synthetic experiments are extensively described in Section 4 and Appendix F and part of the code base at https://anonymous.4open.science/r/perm_equiv_gn_cdes-BBF8/README.md.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All chosen hyperparameters and the ranges we chose them from are detailed in Appendix F.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We repeat each experiment multiple times over different random seeds for data generation and model initialisations. We report either standard deviations or confidence intervals.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide information about the type and amount of compute used for each experiment in Section F.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: As this paper does not involve research involving human subjects or participants, and all data has been properly standardised and anonymised. The work adheres to the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.

- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We briefly discuss the societal impact of this work in Appendix G.1.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper does not pose any of the mentioned risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All datasets, programming languages and libraries used in this work have been attributed properly.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets were introduced in this paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No crowdsourcing nor research with human subjects is part of this paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [NA]

Justification: No LLMs have been used for the core method development of this paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.

A Permutation equivariant linear basis terms

Maron et al. [44] derive the 15 basis elements of all linear permutation equivariant maps $L : \mathbb{R}^{n^2} \rightarrow \mathbb{R}^{n^2}$ on edge-valued signals $\mathbf{A} \in \mathbb{R}^{n^2 \times 1}$. They show the following 15 maps span this space:

We denote by $\mathbf{1} \in \mathbb{R}^n$ the vector whose entries are all equal to one, and by diag the operator that either places a vector on the diagonal of a matrix or extracts the diagonal of a matrix, depending on context. With this notation, the basis elements of $\mathfrak{E}_{\Sigma_n}(2, 1)^1$ can be listed explicitly as follows:

1. The identity and transpose operations: $L(\mathbf{A}) = \mathbf{A}, L(\mathbf{A}) = \mathbf{A}^T$,
2. Elimination of non-diagonal elements: $L(\mathbf{A}) = \text{diag}(\text{diag}(\mathbf{A}))$,
3. Sum of rows replicated on rows/columns/diagonal: $L(\mathbf{A}) = \mathbf{1}\mathbf{1}^T\mathbf{A}, L(\mathbf{A}) = \mathbf{1}(\mathbf{A}\mathbf{1})^T, L(\mathbf{A}) = \text{diag}(\mathbf{A}\mathbf{1})$,
4. Sum of columns replicated on rows/columns/diagonal: $L(\mathbf{A}) = \mathbf{A}^T\mathbf{1}\mathbf{1}^T, L(\mathbf{A}) = \mathbf{1}(\mathbf{A}^T\mathbf{1})^T, L(\mathbf{A}) = \text{diag}(\mathbf{A}^T\mathbf{1})$,
5. Sum of all matrix entries replicated on all entries/diagonal: $L(\mathbf{A}) = \mathbf{1}^T\mathbf{A}\mathbf{1} \cdot \mathbf{1}\mathbf{1}^T, L(\mathbf{A}) = \mathbf{1}^T\mathbf{A}\mathbf{1} \cdot \text{diag}(\mathbf{1})$,
6. Sum of all diagonal entries replicated on all entries/diagonal: $L(\mathbf{A}) = \mathbf{1}^T\text{diag}(\mathbf{A}) \cdot \mathbf{1}\mathbf{1}^T, L(\mathbf{A}) = \mathbf{1}^T\text{diag}(\mathbf{A}) \cdot \text{diag}(\mathbf{1})$,
7. Diagonal elements replicated on rows/columns: $L(\mathbf{A}) = \text{diag}(\mathbf{A})\mathbf{1}\mathbf{1}^T, L(\mathbf{A}) = \mathbf{1}\text{diag}(\mathbf{A})^T$,

B Projection theorem

Definition B.1 (Projection). *Let \mathcal{H} be a Hilbert space with inner product (\cdot, \cdot) . A linear map $p : \mathcal{H} \rightarrow \mathcal{H}$ is called a projection if $p \circ p = p$. A projection is called orthogonal if for all $x \in \mathcal{H}$ it holds that for all $y \in \mathcal{H}$, $p(y) = 0$ implies $(p(x), y) = 0$.*

Intuitively, this is saying that applying a projection operator twice yields the same result as applying it once, which captures the idea of “selecting” or “filtering out” a specific component of a vector. The additional condition for an orthogonal projection ensures that the difference between any vector and its projected counterpart is perpendicular to the subspace, embodying the concept of the closest approximation in the geometry of the space.

Theorem B.2 (Projection Theorem, e.g. [1], Chapter 1). *Let \mathcal{H} be a Hilbert space and U a closed linear subspace. Then every $x \in \mathcal{H}$ can uniquely be written as $x = u + u^*$ where $u \in U$, $u^* \in U^\perp := \{y \in \mathcal{H} : (y, u) = 0 \forall u \in U\}$. Moreover, the map $P_U : \mathcal{H} \rightarrow \mathcal{H}$ defined as $P_U(x) = u$ is an orthogonal projection and satisfies*

$$\|P_U(x) - x\| = \inf_{u' \in U} \|u' - x\|. \quad (9)$$

for all $x \in \mathcal{H}$.

The Projection Theorem asserts that every vector in a Hilbert space can be uniquely decomposed into a sum of two parts: one lying in a closed subspace and the other in its orthogonal complement, thus clarifying the geometric structure of the space. Moreover, it guarantees the projection onto a subspace is the unique map that assigns to each element $x \in \mathcal{H}$ the closest element from that subspace, thereby establishing a rigorous method for optimal approximation in terms of the inner product.

C General theory of equivariance

In the following, we formalise notions of how symmetries can “act” on data and how we want to define functions that respect these symmetries.

C.1 Equivariance of static functions

Definition C.1 (Group Action). *Let G be a group and X a set. A group action of G on X is a function $\cdot : G \times X \rightarrow X$ satisfying:*

1. **Identity:** for all $x \in X$, $e \cdot x = x$, where e is the identity element of G ,
2. **Compatibility:** for all $g, h \in G$ and $x \in X$, $(gh) \cdot x = g \cdot (h \cdot x)$.

This definition formalises the idea that the elements of the group G can act on the set X in a way that reflects the group's structure. It essentially captures the concept of symmetry by showing how each group element transforms or "moves" elements in X .

Definition C.2 (Group Representation). *Let G be a group and V a vector space. A group representation of G on V is a homomorphism $\rho : G \rightarrow \text{GL}(V)$ such that:*

- **Identity:** $\rho(e) = I_V$, where e is the identity element of G and I_V is the identity map on V ,
- **Compatibility:** For all $g_1, g_2 \in G$, $\rho(g_1 g_2) = \rho(g_1) \rho(g_2)$.

This definition realises abstract group elements as concrete, invertible linear transformations on V , thereby representing the group's symmetry in a linear framework. It provides a way to analyse the structure of G by studying how it acts on a vector space.

Definition C.3 (Equivariance). *Let $\rho : G \rightarrow \text{GL}(V)$ and $\sigma : G \rightarrow \text{GL}(W)$ be representations of a group G on vector spaces V and W , respectively. A function $f : V \rightarrow W$ is said to be G -equivariant if for all $g \in G$ and $v \in V$,*

$$f(\rho(g)(v)) = \sigma(g)(f(v)).$$

This definition states that applying the group action before or after the function f produces the same outcome. Equivariance ensures that f preserves the symmetry structure imposed by G , making it a natural and symmetry-respecting mapping between the spaces.

C.2 Haar measure

Definition C.4 (Haar measure). *Let G be a compact Hausdorff topological group. There exists a unique regular Borel measure μ on G invariant under both left and right translations:*

$$\mu(gE) = \mu(Eg) = \mu(E) \quad \text{for all } g \in G, E \subseteq G \text{ Borel}$$

with total mass 1, i.e. $\mu(G) = 1$. We call this measure the Haar measure.

Remark C.5. *If G is finite of order $|G|$, the Haar measure reduces to the uniform distribution:*

$$\mu(\{g\}) = \frac{1}{|G|} \quad \forall g \in G, \quad \int_G f(g) d\mu(g) = \frac{1}{|G|} \sum_{g \in G} f(g).$$

In this discrete (hence compact) setting, Haar integration is exactly averaging over group elements.

C.3 Projection of non-equivariant functions

Given Hilbert spaces V and W , the set of functions $\{f : V \rightarrow W\}$ form a vector space under pointwise scalar multiplication and addition. Given group representations ρ and σ of a group G over V and W , respectively, we can consider the subspace of functions that are equivariant with respect to these representations. The following lemma and corollary show what the projection onto this subspace looks like.

Lemma C.6. *Let V and W be finite-dimensional vector spaces with an action of a finite group G with representations $\rho : G \rightarrow \text{GL}(V)$ and $\sigma : G \rightarrow \text{GL}(W)$ respectively. For any linear map $T : V \rightarrow W$, define the averaged map by*

$$T_{\text{avg}}(v) = \int_G \sigma(g)^{-1} T(\rho(g)v) d\mu(g).$$

Then T_{avg} is G -equivariant, and moreover, T is G -equivariant if and only if $T = T_{\text{avg}}$.

Proof. To show that T_{avg} is G -equivariant, take any $h \in G$ and $v \in V$. Then

$$\begin{aligned} T_{\text{avg}}(\rho(h)v) &= \int_G \sigma(g)^{-1} T(\rho(g)\rho(h)v) d\mu(g) \\ &= \int_G \sigma(g)^{-1} T(\rho(gh)v) d\mu(g). \end{aligned}$$

By substituting $g' = gh$ (so that $g = g'h^{-1}$), we obtain

$$\begin{aligned} T_{\text{avg}}(\rho(h)v) &= \int_G \sigma(g'h^{-1})^{-1} T(\rho(g')v) \, d\mu(g') \\ &= \int_G \sigma(h)\sigma(g')^{-1} T(\rho(g')v) \, d\mu(g') \\ &= \sigma(h)T_{\text{avg}}(v), \end{aligned}$$

where we used $\sigma(g'h^{-1})^{-1} = \sigma(h)\sigma(g')^{-1}$ because σ is a homomorphism. Thus, T_{avg} is G -equivariant.

Next, if T is already G -equivariant, then for every $g \in G$,

$$\sigma(g)^{-1} T(\rho(g)v) = T(v),$$

so that

$$T_{\text{avg}}(v) = \int_G T(v) = T(v).$$

□

Corollary C.7. *In the setting of the previous lemma, the map $(\)_{\text{avg}} : \text{Hom}(V, W) \rightarrow \text{Hom}(V, W)$ defined by $T \mapsto T_{\text{avg}}$ is the projection map onto the subspace of G -equivariant linear maps.*

Proof. In the previous proof we saw that $(\)_{\text{avg}}$ acts as the identity on the subspace of G -equivariant linear maps. Since applying $(\)_{\text{avg}}$ to any linear map T yields a G -equivariant map T_{avg} , $(\)_{\text{avg}}$ is a projection onto the subspace of G -equivariant linear maps. □

Next, we study what form the projection of the composition of a linear equivariant function with another function takes.

Proposition C.8. *Let V, X, W be finite-dimensional vector spaces with an action of a finite group G given by representations $\rho_V : G \rightarrow \text{GL}(V)$, $\rho_X : G \rightarrow \text{GL}(X)$ and $\rho_W : G \rightarrow \text{GL}(W)$ respectively. Denote the projection of $\text{Hom}(V, W)$ onto the space of G -equivariant functions with respect to the representations ρ_V and ρ_W by $P_{V,W}$ and similar for the other combinations. Now suppose we have functions $T : V \rightarrow W$, $R : V \rightarrow X$, $S : X \rightarrow W$ such that*

- $T = S \circ R$,
- S is G -equivariant and linear.

We then have $P_{V,W}(T) = S \circ P_{V,X}(R)$.

Proof. Since $P_{V,W}$ projects onto an equivariant subspace, by the above we know that

$$(P_{V,W}T)(v) = \int_G \rho_W(g)^{-1} T(\rho_V(g)v) \, d\mu(g)$$

for all $v \in V$. Now by definition and the linearity of S ,

$$\begin{aligned} (P_{V,W}T)(v) &= \int_G \rho_W(g)^{-1} (S \circ R)(\rho_V(g)v) \, d\mu(g) \\ &= \int_G \rho_W(g)^{-1} S(R(\rho_V(g)v)) \, d\mu(g) \\ &= \int_G S(\rho_X(g)^{-1} R(\rho_V(g)v)) \, d\mu(g) \quad (\text{equivariance of } S) \\ &= S \left(\int_G \rho_W(g)^{-1} R(\rho_V(g)v) \, d\mu(g) \right) \quad (\text{linearity of } S) \\ &= S((P_{V,X}R)(v)) \end{aligned}$$

since $S(\rho_X(g)^{-1} R(\rho_V(g)v)) = \rho_W(g)^{-1} S(R(\rho_V(g)v))$ by the equivariance of S . This concludes the proof. □

C.4 Equivariance in Neural ODEs

Definition C.9. Let $f : X \times I \rightarrow X$ be a smooth function, where $I \subset \mathbb{R}$ is an interval containing 0 and X is some Banach space. For each $x \in X$, consider the initial value problem

$$\frac{dz}{dt} = f(z, t), \quad z(0) = x. \quad (10)$$

The flow of the ODE is the mapping $\Phi : I \times X \rightarrow X$, $(t, x) \mapsto \Phi_t(x)$, which assigns to each pair (t, x) the solution to the IVP in equation 10 with initial condition x evaluated at time t .

Proposition C.10. Let $\Phi_t(x)$ be the flow of an ODE with vector field $f : X \times I \rightarrow X$ and let P be the projection operator of the space of functions $\{g : X \rightarrow X\}$ onto the space G -equivariant maps with respect to a group representation ρ of G on X . Then the projected flow $P(\Phi_t)$ and the flow $\tilde{\Phi}_t$ of the projected vector field $P(f) : X \times I \rightarrow X$ agree for all times $t \in I$, i.e. $P(\Phi_t) = \tilde{\Phi}_t$.

Proof. First, we note that by Corollary C.7, P can be realised by

$$\begin{aligned} P(f)(x, t) &= \int_G \rho(g)^{-1} f(\rho(g)x, t) \, d\mu(g), \\ P(\Phi_t)(x) &= \int_G \rho(g)^{-1} \Phi_t(\rho(g)x) \, d\mu(g). \end{aligned}$$

Now, by the uniqueness theory of ODEs (see for example Theorem 6 in [43]), it suffices to show that $P(\Phi_t)$ satisfies the same initial condition and ODE as $\tilde{\Phi}_t$. To that end, we calculate

$$\begin{aligned} \frac{d}{dt}(P\Phi_t)(x) &= \int_G \rho(g)^{-1} \frac{d}{dt} \Phi_t(\rho(g)x) \, d\mu(g) \\ &= \int_G \rho(g)^{-1} f(\rho(g)x, t) \, d\mu(g) \\ &= P(f)(x, t) \end{aligned}$$

and

$$\begin{aligned} (P\Phi_0)(x) &= \int_G \rho(g)^{-1} \Phi_0(\rho(g)x) \, d\mu(g) \\ &= \int_G \rho(g)^{-1} \rho(g)x \, d\mu(g) \\ &= x. \end{aligned}$$

This finishes the proof. □

This theorem shows that, under the stated conditions, if one projects the full flow Φ_t using P , then this projected flow is exactly the same as the flow obtained by integrating the projected vector field $P(f)$.

D Application to Graph Neural CDE

Definition D.1 (Representation of the Permutation Group). We define a representation of the permutation group S_n on the space of node features $\mathbb{R}^{n \times d}$ for each $\sigma \in S_n$ by $\rho(\sigma)(\mathbf{X}) = \mathbf{P}_\sigma \mathbf{X}$ where \mathbf{P}_σ is the permutation matrix associated with σ . Similarly, we define a representation on the space of node features $\mathbb{R}^{n \times n}$ by $\pi(\sigma)(\mathbf{X}) = \mathbf{P}_\sigma \mathbf{X} \mathbf{P}_\sigma^T$.

The pre-multiplication by \mathbf{P}_σ permutes the rows of a matrix, and the post-multiplication permutes the columns.

Proof of Theorem 3.1. Firstly, by Proposition C.10, the projection of any Graph Neural ODE onto the subspace of permutation equivariant functions is the Neural ODE where the vector field is the projection of the Graph Neural ODE. Hence, it suffices to consider the projection of the integrand in Equation 4.

We can decompose this integrand as the composition of the following two functions:

$$R(\mathbf{Z}, \mathbf{A}, \frac{d\mathbf{A}}{ds}) = \left(\mathbf{Z}, \mathbf{W}^{(DR)} \left[\begin{array}{c} \mathbf{A} \\ \frac{d\mathbf{A}}{ds} \end{array} \right] \right)$$

$$S(\mathbf{Z}, \tilde{\mathbf{A}}) = \tilde{\mathbf{A}} \mathbf{Z}^{(L)} \mathbf{W}^{(L)}$$

where $\mathbf{Z}^{(L)}$ are latent representations obtained by an iterative convolutional process of the form $\mathbf{Z}^{(l)} = \tilde{\mathbf{A}} \mathbf{Z}^{(l-1)} \mathbf{W}^{(l-1)}$ for $l \in \{1, \dots, L\}$ where the $\mathbf{W}^{(l)}$ are learnable matrices.

The function S is a standard convolutional GNN and hence both permutation equivariant and linear in \mathbf{Z} . Hence, we can apply Proposition C.8 and realise we only need to consider the projection of R .

But since $W^{(DR)}$ is learnable and in Section 2.3 we characterised the space of linear equivariant functions on the space of adjacency matrices, instead of projecting onto the space of equivariant functions, we may directly parameterise an element of this space as a linear combination of the basis elements. This concludes the proof. \square

E Proofs

Proof the GN-CDE (Equation 3) is not permutation equivariant. For notational simplicity, we flatten $f_\theta(\mathbf{Z}_s, \mathbf{A}_s)$ and the control $d\hat{\mathbf{A}}_s$ into tensors of ranks 2 and 1, respectively:

$$F_\theta(\mathbf{Z}_s, \mathbf{A}_s) = \text{vec}[f_\theta(\mathbf{Z}_s, \mathbf{A}_s)] \in \mathbb{R}^{nd_z \times 2n^2}, \quad (11)$$

$$\mathbf{B} = \text{vec}[d\hat{\mathbf{A}}_s] \in \mathbb{R}^{2n^2}. \quad (12)$$

Under a node permutation $\sigma \in S_n$, with permutation matrix \mathbf{P}_σ , these transform as

$$F_\theta(\mathbf{Z}_s, \mathbf{A}_s) \mapsto (\mathbf{P}_\sigma \otimes \mathbf{I}_{d_z}) F_\theta(\mathbf{Z}_s, \mathbf{A}_s), \quad (13)$$

$$\mathbf{B} \mapsto (\mathbf{P}_\sigma \otimes \mathbf{P}_\sigma \otimes \mathbf{I}_2) \mathbf{B}. \quad (14)$$

Requiring that the contraction

$$T_\theta(\mathbf{Z}_s, \mathbf{A}_s) = F_\theta(\mathbf{Z}_s, \mathbf{A}_s) \mathbf{B}$$

be permutation-equivariant then gives

$$T_\theta(\mathbf{P}_\sigma \mathbf{Z}_s, \mathbf{P}_\sigma \mathbf{A}_s \mathbf{P}_\sigma^\top) = (\mathbf{P}_\sigma \otimes \mathbf{I}_{d_z}) T_\theta(\mathbf{Z}_s, \mathbf{A}_s) (\mathbf{P}_\sigma^\top \otimes \mathbf{P}_\sigma^\top \otimes \mathbf{I}_2) \quad \forall \sigma \in S_n.$$

No general architecture for f_θ (for example, a standard multilayer perceptron) can satisfy this constraint unless it is explicitly parametrised to output elements in the space of permutation-equivariant linear maps. \square

Proof the GN-CDE fusion operation (Equation 4) is not permutation equivariant. Firstly we note that by writing $\mathbf{W}^{(F)} = \begin{bmatrix} \mathbf{W}_1^{(F)} & \mathbf{W}_2^{(F)} \end{bmatrix}$ for $\mathbf{W}_1^{(F)}, \mathbf{W}_2^{(F)} \in \mathbb{R}^{n \times n}$, the fusion takes the form

$$\tilde{\mathbf{A}}_s = \mathbf{W}^{(F)} \begin{bmatrix} \mathbf{A}_s \\ \frac{d\mathbf{A}_s}{ds} \end{bmatrix} = \begin{bmatrix} \mathbf{W}_1^{(F)} & \mathbf{W}_2^{(F)} \end{bmatrix} \begin{bmatrix} \mathbf{A}_s \\ \frac{d\mathbf{A}_s}{ds} \end{bmatrix} = \mathbf{W}_1^{(F)} \mathbf{A}_s + \mathbf{W}_2^{(F)} \frac{d\mathbf{A}_s}{ds}. \quad (15)$$

Thus, the fusion operation can be viewed as pre-multiplying each of \mathbf{A}_s and $\frac{d\mathbf{A}_s}{ds}$ by a learnable weight matrix and then summing the result. However, letting \mathbf{P} be an arbitrary permutation matrix, we have

$$\mathbf{P} \tilde{\mathbf{A}}_s \mathbf{P}^T = \mathbf{P} \mathbf{W}_1^{(F)} \mathbf{A}_s \mathbf{P}^T + \mathbf{P} \mathbf{W}_2^{(F)} \frac{d\mathbf{A}_s}{ds} \mathbf{P}^T \quad (16a)$$

$$\neq \mathbf{W}_1^{(F)} \mathbf{P} \mathbf{A}_s \mathbf{P}^T + \mathbf{W}_1^{(F)} \mathbf{P} \frac{d\mathbf{A}_s}{ds} \mathbf{P}^T \quad (16b)$$

in general, since $\mathbf{W}_1^{(F)}$ and $\mathbf{W}_2^{(F)}$ do not necessarily commute with \mathbf{P} . \square

Proof of Proposition 3.2. First we show permutation equivariance which follows immediately from the construction of the linear fusion. For let \mathbf{P} be an arbitrary permutation matrix. Then

$$L_1(\mathbf{P}\mathbf{A}_s\mathbf{P}^T) + L_2\left(\mathbf{P}\frac{d\mathbf{A}_s}{ds}\mathbf{P}^T\right) = \mathbf{P}L_1(\mathbf{A}_s)\mathbf{P}^T + \mathbf{P}L_2\left(\frac{d\mathbf{A}_s}{ds}\right)\mathbf{P}^T = \mathbf{P}\bar{\mathbf{A}}_s\mathbf{P}^T. \quad (17)$$

Since ζ_θ is implemented as a GNN, $\zeta_\theta(\mathbf{P}\mathbf{X}_{t_0}, \mathbf{P}\mathbf{A}_{t_0}\mathbf{P}^T) = \mathbf{P}\zeta_\theta(\mathbf{X}_{t_0}, \mathbf{A}_{t_0}) = \mathbf{P}\mathbf{Z}_{t_0}$ and so overall for all $t \in (t_0, t_N]$ we have

$$\mathbf{P}\mathbf{Z}_{t_0} + \int_{t_0}^t \sigma(\mathbf{P}\bar{\mathbf{A}}_s\mathbf{P}^T\mathbf{P}\mathbf{Z}_s^{(L)}\mathbf{W}^{(L)})ds = \mathbf{P}\mathbf{Z}_{t_0} + \int_{t_0}^t \mathbf{P}\sigma(\bar{\mathbf{A}}_s\mathbf{Z}_s^{(L)}\mathbf{W}^{(L)})ds \quad (18a)$$

$$= \mathbf{P}\mathbf{Z}_{t_0} + \mathbf{P} \int_{t_0}^t \sigma(\bar{\mathbf{A}}_s\mathbf{Z}_s^{(L)}\mathbf{W}^{(L)})ds \quad (18b)$$

$$= \mathbf{P} \left(\mathbf{Z}_{t_0} + \int_{t_0}^t \sigma(\bar{\mathbf{A}}_s\mathbf{Z}_s^{(L)}\mathbf{W}^{(L)})ds \right) \quad (18c)$$

$$= \mathbf{P}\mathbf{Z}_t \quad (18d)$$

using the orthogonality of \mathbf{P} and linearity of the integral.

To show time-warp equivariance, suppose a latent path $\mathbf{Z} : [0, T] \rightarrow \mathbb{R}^{n \times d_z}$ satisfies

$$\mathbf{Z}(0) = \mathbf{Z}_0, \quad \frac{d\mathbf{Z}}{dt}(t) = f_\theta(\mathbf{Z}(t), \mathbf{A}(t)) \frac{d\hat{\mathbf{A}}}{dt}(t). \quad (19)$$

for some dynamic graph data $\mathbf{A} : [0, T] \rightarrow \mathbb{R}^{n \times n}$. Then the warped path $\hat{\mathbf{Z}} := \mathbf{Z} \circ \tau$ satisfies

$$\hat{\mathbf{Z}}(0) = (\mathbf{Z} \circ \tau)(0) = \mathbf{Z}(0) = \mathbf{Z}_0. \quad (20)$$

and by a simple application of the chain rule we get

$$\frac{d\hat{\mathbf{Z}}}{dt}(t) = \frac{d}{dt}(\mathbf{Z} \circ \tau)(t) = \frac{d\mathbf{Z}}{dt}(\tau(t))\tau'(t) \quad (21a)$$

$$= f_\theta(\mathbf{Z}(\tau(t)), \mathbf{A}(\tau(t))) \frac{d\hat{\mathbf{A}}}{dt}(\tau(t))\tau'(t) \quad (21b)$$

$$= f_\theta(\hat{\mathbf{Z}}(t), \mathbf{A}(\tau(t))) \frac{d}{dt}(\hat{\mathbf{A}} \circ \tau)(t) \quad (21c)$$

which is what we wanted. \square

Proof of equivariance in Section 3.3. First we show that for general $\mathbf{A} \in \mathbb{R}^{n \times m \times d}$ and $\mathbf{B} \in \mathbb{R}^{n \times d}$ we have $\mathbf{P}(\mathbf{A} \odot \mathbf{B}) = \mathbf{P}\mathbf{A} \odot \mathbf{P}\mathbf{B}$. In Einstein notation we have

$$\{\mathbf{P}(\mathbf{A} \odot \mathbf{B})\}_{ij} = \mathbf{P}_{ik}(\mathbf{A} \odot \mathbf{B})_{kj} \quad (22a)$$

$$= \sum_k \mathbf{P}_{ik}\mathbf{A}_{kjl}\mathbf{B}_{kl}. \quad (22b)$$

At the same time

$$(\mathbf{P}\mathbf{A} \odot \mathbf{P}\mathbf{B})_{ij} = (\mathbf{P}\mathbf{A})_{ijl}(\mathbf{P}\mathbf{B})_{il} \quad (23a)$$

$$= \mathbf{P}_{ik}\mathbf{A}_{kjl}\mathbf{P}_{ik'}\mathbf{B}_{k'l} \quad (23b)$$

$$= \sum_k \mathbf{P}_{ik}\mathbf{A}_{kjl}\mathbf{B}_{kl} \quad (23c)$$

$$= \{\mathbf{P}(\mathbf{A} \odot \mathbf{B})\}_{ij} \quad (23d)$$

because $\mathbf{P}_{ik}, \mathbf{P}_{ik'} \in \{0, 1\}$, so the only terms contributing to the sum are those where $\mathbf{P}_{ik} = \mathbf{P}_{ik'}$ which happens if and only if $k = k'$ as \mathbf{P} is a permutation matrix.

Hence, we see that the vector field in equation 7 is permutation equivariant. By the above, we can conclude that the entire model is equivariant. \square

F Implementation details

We will release all code alongside the camera-ready version of this paper. The implementation is written in the Python programming language [61], and uses the JAX framework [4]. Key dependencies include Diffrax [34] for differentiable ODE solvers, Equinox [35] for neural network modules in JAX, Optax [4] for optimisation, and Lineax [49] for linear algebra routines. Additional dependencies include NumPy [28], Exca [50], PyTorch [54], and PyTorch Geometric [23].

All differential equation-based models use the Tsitouras’ 5/4 method [60] as implemented in Diffrax.

F.1 Compute resources

All experiments were conducted on a computing cluster equipped with NVIDIA H100 (80 GB) and H200 (141 GB) GPUs. Each node in the cluster features 64-core AMD EPYC 9334 CPUs running at 3.90 GHz, along with 256 GB of RAM. The heat diffusion and gene regulation experiments described in Section 4.1 required approximately 3.5 GPU days to complete, whereas the oversampling and irregularity experiments finished within a total of 12 GPU hours. The PyTorch Geometric Temporal experiments ran for approximately 8 hours. Running STG-NCDE, GN-CDE, and both versions of our model on the two Temporal Graph Benchmark tasks required around 8 GPU hours in total.

F.2 Heat diffusion and gene regulation

Two important examples of dynamical systems on networks are heat diffusion and gene regulation dynamics, both of which describe the local exchange of information or state between connected nodes. While heat diffusion models symmetric spreading (e.g., temperature equalisation), gene regulation introduces asymmetric, nonlinear interactions typical of biological systems. These processes can be modelled as follows:

$$\text{Heat Diffusion: } \frac{d\mathbf{x}_u(t)}{dt} = \mathcal{L}_t \mathbf{x}_u(t) = \sum_{v \in \mathcal{N}_u(t)} \left(\frac{\mathbf{x}_u(t)}{\sqrt{d_u}} - \frac{\mathbf{x}_v(t)}{\sqrt{d_v}} \right) \quad (24)$$

$$\text{Gene Regulation: } \frac{d\mathbf{x}_u(t)}{dt} = -\mathbf{x}_u(t)f + \sum_{v \in \mathcal{N}_u(t)} \frac{\mathbf{x}_v(t)}{\mathbf{x}_v(t) + 1} \quad (25)$$

where $\mathbf{x}_u(t)$ is the temperature of node u at time t , $\mathcal{N}_u(t)$ denotes its neighbourhood, and \mathcal{L}_t is the normalised graph Laplacian.

F.3 Synthetic experiments

To ensure a fair comparison, the vector fields of all differential equation-based models are implemented using a two-layer GCN [38] with a hidden dimension of 16, and layer normalisation is applied. The models are trained for 2000 epochs using the Adam optimiser [37] with a learning rate of 10^{-2} and weight decay of 10^{-4} .

F.3.1 SIR model

The Susceptible-Infected-Recovered (SIR) model on graphs [33] is a network-based extension of the foundational compartmental model in epidemiology [32] that describes how infectious diseases spread through a structured population. It partitions nodes in a graph into three groups: susceptible (S_v), infected (I_v), and recovered (R_v) for each node v , with transitions governed by the system

$$\frac{dS_v}{dt} = -\beta S_v \sum_{u \in \mathcal{N}(v)} I_u, \quad \frac{dI_v}{dt} = \beta S_v \sum_{u \in \mathcal{N}(v)} I_u - \gamma I_v, \quad \frac{dR_v}{dt} = \gamma I_v, \quad (26)$$

where β is the transmission rate, γ the recovery rate, and $\mathcal{N}(v)$ denotes the neighbours of node v in the graph.

A key insight from the graph-based SIR model is that the epidemic threshold depends not only on the infection ratio $R = \beta/\gamma$, but also on the structure of the underlying graph, such as its spectral

radius or degree distribution. If local connectivity supports sufficient transmission, the infection can spread widely; otherwise, it dies out. Despite its structural complexity, the graph SIR model preserves essential features of epidemic dynamics and extends the classical model to networked populations.

F.3.2 Oversampling

To study the dependence of performance and compute time on the number of observed samples, we construct grid graphs with 100 nodes and generate irregularly sampled time series between $T = 0$ and $T = 1$, using $n = 20, 40, 60, \dots, 200$ observations respectively. The dynamic graph topology is generated as described in Section 4.1. Trajectories are produced by numerically solving Equation 26 with random initial conditions. We use two parameter settings: $(\beta_1 = 0.25, \gamma_1 = 0.7)$, modelling a scenario where the epidemic dies out, and $(\beta_2 = 0.3, \gamma_2 = 0.3)$, modelling a scenario where the infection spreads. For each setting, we generate 50 trajectories each for training, validation, and test sets.

The task is to perform binary classification, predicting whether a given trajectory corresponds to an outbreak (i.e., sustained spread of infection) or a non-outbreak (i.e., rapid die-out).

We train the models using binary cross-entropy loss for 2000 epochs, applying early stopping with a patience of 200 epochs and a minimum of 200 epochs. Training is performed using the Adam optimizer [37] with a learning rate of 10^{-2} and a weight decay of 10^{-4} . The model hyperparameters are summarised in Table 4.

Table 4: Hyperparameters for SIR experiments.

	Perm Equiv GN-CDE	DCRNN
Hidden Dimension	32	8
Number of Layers	3	3
Data Embedding Dimension	3	3
Layer type	GCN	–
Chebyshev Order	–	3

F.3.3 Irregular observation spacing

To assess how the regularity of observation times affects model performance, we generate classification data exactly as before, except that we now sample the irregular observation time-stamps from a Gamma-driven process. In more detail, let N be the total number of desired time-points on the interval $[0, T]$. We begin by drawing $N - 1$ independent inter-arrival increments

$$\Delta t_i \sim \text{Gamma}(k, \theta), \quad i = 1, \dots, N - 1,$$

where $k > 0$ controls the shape (and hence burstiness) and $\theta > 0$ scales the raw increments. To force the observations to span exactly $[0, T]$, we normalise the increments by their sum:

$$\tilde{\Delta} t_i = \frac{\Delta t_i}{\sum_{j=1}^{N-1} \Delta t_j} T, \quad t_i = \sum_{j=1}^i \tilde{\Delta} t_j,$$

with $t_0 = 0$ and $t_N = T$. This construction guarantees $0 = t_0 < t_1 < \dots < t_N = T$ and yields exactly $N + 1$ strictly increasing time-stamps.

In our experiments, we set $T = 1$ and draw $N = 20$ observations per trajectory using this Gamma-based sampler for $k \in \{3, 5, 10, 25, 50, 100\}$. All other aspects of model training - network architecture, optimisation schedule, and hyperparameters - remain identical to those in the previous section.

F.4 Real-world experiments

All models are trained over 200 epochs with an early stopping criterion with a patience of 15 and a minimum epoch of 20. In all applications, we utilise cubic splines [34, Chapter 3.5] for interpolation of the data. All hyperparameters have been selected from the ranges in Table 5 using a grid search. For the chosen values see Tables 6, 7 and 8. All experiments have been averaged over 10 random seeds.

Table 5: Ranges the optimal hyperparameters have been chose from for (Perm Equiv) GN-CDE (top) and STG-NCDE (bottom).

	Range
Hidden Dimension	{8, 16, 32, 64}
Number of Layers	{2, 3, 4}
Data Embedding Dimension	{8, 16, 32, 64}
Number of Recurrent Layers	{1, 2, 3}
Recurrent Hidden Dimension	{8, 16, 32, 64}
Number of GNN Layers	{1, 2, 3}
GNN Hidden Dimension	{8, 16, 32, 64}
Data Embedding Dimension	{8, 16, 32, 64}
Chebychev Order	{2}
Node Embedding Dimension	{5}

Table 6: Hyperparameters for real-world experiments for the GN-CDE model.

	tgbn-trade	tgbn-genre	twitter-tennis	england-covid
Hidden Dimension	8	8	8	32
Number of Layers	3	2	2	4
Data Embedding Dimension	16	8	-	-
Layer type	GCN	GCN	GCN	GCN

F.5 Source target identification

In [58], the authors address the well-known issue that machine-learned models underperform on node affinity prediction tasks compared to simple heuristics [52] by modifying the message construction in TGN [52] as follows: Given an interaction between nodes u and v at time t with associated data $e_{uv}(t)$, the messages $\mathbf{m}_u(t)$ and $\mathbf{m}_v(t)$ sent to nodes u and v , respectively, are defined as

$$\mathbf{m}_u(t) = \text{msg}_s(\mathbf{s}_u(t^-), \mathbf{s}_v(t^-), \phi_t(\Delta t), e_{uv}(t), \phi_n(u), \phi_n(v)), \quad (27)$$

$$\mathbf{m}_v(t) = \text{msg}_d(\mathbf{s}_v(t^-), \mathbf{s}_u(t^-), \phi_t(\Delta t), e_{uv}(t), \phi_n(v), \phi_n(u)), \quad (28)$$

where Δt is the time elapsed since their last interaction and $\phi : \mathbb{R} \rightarrow \mathbb{R}^n$ is an encoder function for node indices. In their work, the authors use $\phi(u) = (\cos(\omega_i u))_{i=0}^{n-1}$.

Inspired by this, we incorporate source-target identification into our equivariant fusion by transforming the adjacency matrix \mathbf{A}_s and its derivative $\frac{d\mathbf{A}_s}{ds}$ using an attention-style approach. Specifically, each entry a_{uv} of either matrix is passed through an MLP together with the encoding of both u and v to form a new matrix \mathbf{B} with entries

$$b_{uv} = \text{MLP}(a_{uv} || \phi(u) || \phi(v)). \quad (29)$$

where $||$ denotes concatenation. We selected the following hyperparameters: an embedding dimension of $n = 512$, an MLP width of 8, and 2 hidden layers. We slightly deviated from the original approach by defining $\phi(u) = \text{MLP}(u)$ with a hidden dimension of 8 and 2 hidden layers.

Table 7: Hyperparameters for real-world experiments for the Permutation Equivariant GN-CDE model.

	tgbn-trade	tgbn-genre	twitter-tennis	england-covid
Hidden Dimension	32	16	64	64
Number of Layers	2	3	3	3
Data Embedding Dimension	8	8	-	-
Layer type	GCN	GCN	GCN	GCN

Table 8: Hyperparameters for real-world experiments for the STG-NCDE model.

	tgbn-trade	tgbn-genre	twitter-tennis	england-covid
Number of Recurrent Layers	1	3	3	1
Recurrent Hidden Dimension	8	16	32	32
Number of GNN Layers	2	1	2	3
GNN Hidden Dimension	8	16	32	32
Data Embedding Dimension	32	8	-	-
Chebyshev Order	2	2	2	2
Node Embedding Dimension	5	5	5	5

Table 9: Fusion operation weights of the Permutation Equivariant GN-CDE model. Weights with absolute value larger than 0.1 are in bold.

Operation	Layer 1		Layer 2	
	A_s	dA_s	A_s	dA_s
Identity	1.6129	0.9591	1.0747	0.9100
Transpose	0.1026	-0.0210	-0.0270	0.0535
Diagonalisation	0.0297	-0.0754	0.0072	0.0317
Sum rows on	rows	0.7753	-0.0635	0.0358
	columns	0.3407	0.0015	-0.0242
	diagonal	-0.0875	-0.0222	0.0790
Sum all on	all	0.4315	0.4260	0.0039
	diagonal	-0.0002	-0.0001	0.0020
			0.0362	0.0088

G Miscellaneous

G.1 Societal impact

This work constitutes foundational research and is not tied to any specific application or deployment. As such, it shares the general risks and benefits inherent to novel machine-learning architectures. For example, dynamic graph representation learning models have traditionally been applied to tasks such as traffic forecasting [11, 10] and molecular modelling [57], and we hope that the advances presented here will further drive progress in these and related areas. Overall, we assess the societal impact of this work to be predominantly positive.

H Ablation studies

To understand how our model learns to fuse the adjacency matrix and its derivative, we analyse the learned weights for the 15 basis terms (Table 9) on the heat diffusion task. Comparing the two GCN layers, layer 1 assigns importance to both identity and non-identity terms, while layer 2 focuses primarily on the identity operation (favouring the adjacency matrix slightly over its derivative). In layer 1, significant weights are assigned to summing the adjacency matrix rows and replicating them across the entire matrix, suggesting a focus on capturing global graph properties. Similarly, large weights are given to the overall sums of both the adjacency matrix and its derivative.