
Capturing Graphs with Hypo-Elliptic Diffusions

Csaba Toth*
Mathematical Institute
University of Oxford
toth@maths.ox.ac.uk

Darrick Lee*
Mathematical Institute
University of Oxford
leed@maths.ox.ac.uk

Celia Hacker
Department of Mathematics
EPFL
celia.hacker@epfl.ch

Harald Oberhauser
Mathematical Institute
University of Oxford
oberhauser@maths.ox.ac.uk

Abstract

Convolutional layers within graph neural networks operate by aggregating information about local neighbourhood structures; one common way to encode such substructures is through random walks. The distribution of these random walks evolves according to a diffusion equation defined using the graph Laplacian. We extend this approach by leveraging classic mathematical results about hypo-elliptic diffusions. This results in a novel tensor-valued graph operator, which we call the hypo-elliptic graph Laplacian. We provide theoretical guarantees and efficient low-rank approximation algorithms. In particular, this gives a structured approach to capture long-range dependencies on graphs that is robust to pooling. Besides the attractive theoretical properties, our experiments show that this method competes with graph transformers on datasets requiring long-range reasoning but scales only linearly in the number of edges as opposed to quadratically in nodes.

1 Introduction

Obtaining a latent description of the non-Euclidean structure of a graph is central to many applications. One common approach is to construct a set of features for each node that represents the local neighborhood of this node; pooling these node features then provides a latent description of the whole graph. A classic way to arrive at such node features is by random walks: at the given node one starts a random walk, and extracts a summary of the local neighbourhood from its sample trajectories. We revisit this random walk construction and are inspired by two classical mathematical results:

Hypo-elliptic Laplacian. In the Euclidean case of Brownian motion $B = (B_t)_{t \geq 0}$ evolving in \mathbb{R}^n , the quantity $u(t, x) = \mathbb{E}[f(B_t) | B_0 = x]$ solves the heat equation $\partial_t u = \Delta u$ on $[0, \infty) \times \mathbb{R}^n$, $u(0, x) = f(x)$. Seminal work of Gaveau [23] in the 1970s shows that if one replaces $f(B_t)$ in the expectation by a functional of the whole trajectory, $F(B_s : s \in [0, t])$, then a path-dependent heat equation can be derived where the classical Laplacian Δ is replaced by the hypo-elliptic Laplacian.

Free Algebras. A simple way to capture a sequence – for us, a sequence of nodes visited by a random walk on a graph – is to associate with each sequence element an element in an algebra¹ and multiply these algebra elements together. If the algebra multiplication is commutative,

*Equal contribution; order determined by random coin flip.

¹An algebra is a vector space where one can multiply elements; e.g. the set of $n \times n$ matrices with matrix multiplication. This multiplication can be non-commutative; e.g. $A \cdot B \neq B \cdot A$ for general matrices A, B .

the sequential structure is lost but if it is non-commutative, this captures the order in the sequence. In fact, by using the free associative algebra, this can be done faithfully and linear functionals of this algebra correspond to functionals on the space of sequences.

We leverage these ideas from the Euclidean case of \mathbb{R}^d to the non-Euclidean case of graphs. In particular, we construct features for a given node by sampling from a random walk started at this node, but instead of averaging over the end points, we average over path-dependent functions. Put informally, instead of asking a random walker that started at a node, "*What do you see now?*" after k steps, we ask "*What have you seen along your way?*". The above notions from mathematics about the hypo-elliptic Laplacian and the free algebra allow us to formalize this in the form of a generalized graph diffusion equation and we develop algorithms that make this a scalable method.

Related Work. From the ML literature, [49, 27] popularized the combination of deep learning architectures to capture random walk histories. Such ideas have been incorporated, sometimes implicitly, into graph neural networks (GNN) [53, 9, 54, 18, 29, 5, 33] that in turn build on convolutional approaches [37, 38, 27], as well as their combination with attention or message passing [45, 65, 24], and more recent improvements [72, 46, 43, 15] that provide and improve on theoretical guarantees. Another classic approach are graph kernels, see [7] for a recent survey; in particular, the seminal paper [36] explored the connection between diffusion equations and random walkers in a kernel learning context. More recently, [14] proposed sequence kernels to capture the random walk history. Furthermore, [17] uses the signature kernel maximum mean discrepancy (MMD) [16] as a metric for trees which implicitly relies on the algebra of tensors that we use, and [48] aggregates random walk histories to derive a kernel for graphs. Moreover, the concept of network motifs [44, 55] relates to similar ideas that describe a graph by node sequences. Further, the Bethe Hessian [52] has been successfully used in spectral clustering and shares the same goal of capturing pathdependence via a "deformed Laplacians", although the mathematical approach is very different to ours. Directly related to our approach is the topic of learning diffusion models [35, 13, 60, 20, 11] on graphs. While similar ideas on random walks and diffusion for graph learning have been developed by different communities, our proposed method leverages these perspectives by capturing random walk histories through a novel diffusion operation.

Our main mathematical influence is the seminal work of Gaveau [23] from the 1970s that shows how Brownian motion can be lifted into a Markov process evolving in a free algebra to capture path-dependence. This leads to a heat equation governed by the hypo-elliptic Laplacian. These insights had a large influence in PDE theory, see [51, 31], but it seems that their discrete counterpart on graphs has received no attention despite the well-developed literature on random walks on graphs and general non-Euclidean objects, [68, 19, 26, 63]. A key challenge to go from theory to application is handling the computational complexity. To do so, we build on ideas from [62] to design effective algorithms for the hypo-elliptic graph diffusion.

Contribution and Outline. We introduce the hypo-elliptic graph Laplacian which allows to effectively capture random walk histories through a generalized diffusion model.

- In Section 3, we introduce the hypo-elliptic variants of standard graph matrices such as the adjacency matrix and (normalized) graph Laplacians. These hypo-elliptic variants are formulated in terms of tensor-valued matrices rather than scalar-valued matrices, and can be manipulated using linear algebra in the same manner as the classical setting.
- The hypo-elliptic Laplacian leads to a corresponding diffusion model, and in Theorem 1, we show that the solution to this generalized diffusion equation summarizes the microscopic picture of *the entire history of random walks* and not just their location after k steps.
- This solution provides a rich description of the local neighbourhood about a node, which can either be used directly as node features or be pooled over the graph to obtain a latent description of the graph. Theorem 2 shows that these node features characterize random walks on the graph, and we provide an analogous statement for graph features in Appendix E.
- One can solve the hypo-elliptic graph diffusion equation directly with linear algebra, but this is computationally prohibitive and Theorem 3 provides an efficient low-rank approximation.
- Finally, Section 5 provides experiments and benchmarks. A particular focus is to test the ability of our model to capture long-range interactions between nodes and the robustness of pooling operations which makes it less susceptible to the "over-squashing" phenomenon [2].

2 Sequence Features by Non-Commutative Multiplication.

We define the *space of sequences in \mathbb{R}^d* by

$$\text{Seq}(\mathbb{R}^d) := \bigcup_{k=0}^{\infty} (\mathbb{R}^d)^{k+1},$$

where elements are sequences denoted by $\mathbf{x} = (x_0, x_1, \dots, x_k) \in (\mathbb{R}^d)^{k+1}$. Assume we are given an injective map, which we call the *algebra lifting*,

$$\varphi : \mathbb{R}^d \rightarrow H.$$

from \mathbb{R}^d into an algebra H . We can use this to define a *sequence feature map*²

$$\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H, \quad \tilde{\varphi}(\mathbf{x}) = \varphi(\delta_0 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}), \quad (1)$$

where $\delta_0 \mathbf{x} = x_0$ and $\delta_i \mathbf{x} := x_i - x_{i-1}$ for $i \geq 1$ are used to denote the *increments* of a sequence $\mathbf{x} = (x_0, \dots, x_k)$. This map associates to any sequence $\mathbf{x} \in \text{Seq}(\mathbb{R}^d)$ an element of the algebra H . If the multiplication in H is commutative, then the map $\tilde{\varphi}$ would have no information about the order of increments, i.e. $\varphi(\delta_0 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}) = \varphi(\delta_{\pi(0)} \mathbf{x}) \cdots \varphi(\delta_{\pi(k)} \mathbf{x})$ for any permutation π of $\{0, \dots, k\}$. However, if the multiplication in H is "non-commutative enough" we expect $\tilde{\varphi}$ to be injective.

A Free Construction. Many choices for H are possible, but intuitively it makes sense to use the "most general object" for H . The mathematically rigorous approach is to use the *free algebra over \mathbb{R}^d* and we give a summary in Appendix B. Despite this abstract motivation, the algebra H has a concrete form: it is realized as a sequence of tensors in \mathbb{R}^d of increasing degree, and is defined by

$$H := \{\mathbf{v} = (\mathbf{v}_0, \mathbf{v}_1, \mathbf{v}_2, \dots) : \mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}, m \in \mathbb{N}, \|\mathbf{v}\| < \infty\},$$

where by convention $(\mathbb{R}^d)^{\otimes 0} = \mathbb{R}$, and we describe the norm $\|\mathbf{v}\|$ in the paragraph below. For example, if $\mathbf{v} = (\mathbf{v}_m)_{m \geq 0} \in H$, then \mathbf{v}_0 is a scalar, \mathbf{v}_1 is a vector, $\mathbf{v}_2 \in (\mathbb{R}^d)^{\otimes 2}$ is a $d \times d$ matrix, and so on. The vector space structure of H is given by addition and scalar multiplication according to

$$\mathbf{v} + \mathbf{w} := (\mathbf{v}_m + \mathbf{w}_m)_{m \geq 0} \in H \quad \text{and} \quad \lambda \mathbf{v} := (\lambda \mathbf{v}_m)_{m \geq 0} \in H$$

for $\lambda \in \mathbb{R}$, and the algebra structure is given by

$$\mathbf{v} \cdot \mathbf{w} := \left(\sum_{i=0}^m \mathbf{v}_i \otimes \mathbf{w}_{m-i} \right)_{m \geq 0} \in H. \quad (2)$$

An Inner Product. If e^1, \dots, e^d is a basis of \mathbb{R}^d , then every tensor $\mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}$ can be written as

$$\mathbf{v}_m = \sum_{1 \leq i_1, \dots, i_m \leq d} c_{i_1, \dots, i_m} e^{i_1} \otimes \dots \otimes e^{i_m}.$$

This allows us to define an inner product $\langle \cdot, \cdot \rangle_m$ on $(\mathbb{R}^d)^{\otimes m}$ by extending

$$\langle e^{i_1} \otimes \dots \otimes e^{i_m}, e^{j_1} \otimes \dots \otimes e^{j_m} \rangle_m = \begin{cases} 1 & : i_1 = j_1, \dots, i_m = j_m, \\ 0 & : \text{otherwise.} \end{cases} \quad (3)$$

to $(\mathbb{R}^d)^{\otimes m}$ by linearity. This gives us an inner product on H ,

$$\langle \mathbf{v}, \mathbf{w} \rangle := \sum_{m \geq 0} \langle \mathbf{v}_m, \mathbf{w}_m \rangle_m$$

such that H is a Hilbert space; in particular we get a norm $\|\mathbf{v}\| := \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$. To sum up, the space H has a rich structure: it has a vector space structure, it has an algebra structure (a noncommutative product), and it is a Hilbert space (an inner product between elements of H gives a scalar).

²There are variants of this sequence feature map, which are discussed in Appendix G.

Characterizing Random Walks. From Equation (1), we have constructed a map $\tilde{\varphi}$ that maps a sequence $\mathbf{x} \in \text{Seq}(\mathbb{R}^d)$ of arbitrary length into the space H (see Appendix C for further details). Our aim is to apply this to the sequence of node attributes corresponding to random walks on a graph. Therefore, the expectation of $\tilde{\varphi}$ should be able to characterize the distribution of the random walk. Formally the map $\tilde{\varphi}$ is *characteristic* if the map $\mu \mapsto \mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})]$ from the space of probability measures on $\text{Seq}(\mathbb{R}^d)$ into H is injective. Indeed, if the chosen lifting φ satisfies some mild conditions this holds for $\tilde{\varphi}$; see Appendix C and [16, 62].

Linear Functionals. The quantity $\mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})]$ characterizes the probability measure μ but is valued in the infinite-dimensional Hilbert space H . Using the inner product, we can instead consider

$$\langle \ell, \mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})] \rangle \text{ for } \ell = (\ell_0, \ell_1, \ell_2, \dots, \ell_M, 0, \dots) \in H \text{ and } M \geq 1 \quad (4)$$

which is equivalent to knowing $\mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})]$; i.e. the set (4) characterizes μ . This is analogous to how one can use either the moment generating function of a real-valued random variable or its sequence of moments to characterize its distribution; the former is one infinite-dimensional object (a function), the latter is a infinite sequence of scalars. We extend a key insight from [62] in Section 4: a linear functional $\langle \ell, \mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})] \rangle$ can be efficiently approximated without directly computing $\mathbb{E}_{\mathbf{x} \sim \mu}[\tilde{\varphi}(\mathbf{x})]$ or storing large tensors.

The Tensor Exponential. While we will continue to keep φ arbitrary for our main results (see [62] and Appendix G for other choices), we will use the *tensor exponential* $\text{exp}_{\otimes} : \mathbb{R}^d \rightarrow H$, defined by

$$\text{exp}_{\otimes}(x) = \left(\frac{x^{\otimes m}}{m!} \right)_{m \geq 0}, \quad (5)$$

as the primary example throughout this paper and in the experiments in Section 5. With this choice, the induced sequence feature map is the discretized version of a classical object in analysis, called the path signature, see Appendix C.

3 Hypo-Elliptic Diffusions

Throughout this section, we fix a labelled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, f)$, that is \mathcal{V} is a set of n nodes $\mathcal{V} = \{1, \dots, n\}$, \mathcal{E} denotes edges and $f : \mathcal{V} \rightarrow \mathbb{R}^d$ is the set of continuous node attributes³ which map each node to an element in the vector space \mathbb{R}^d . Two nodes $i, j \in \mathcal{V}$ are *adjacent* if $(i, j) \in \mathcal{E}$ is an edge, and we denote this by $i \sim j$. The *adjacency matrix* A of a graph is defined by $A_{i,j} = 1$, whenever $i \sim j$, and 0 otherwise. We denote by $\text{deg}(i)$ the number nodes that are adjacent to node i .

Random Walks on Graphs. Let $(B_k)_{k \geq 0}$ be the simple random walk on the nodes \mathcal{V} of \mathcal{G} , where the initial node is chosen uniformly at random. The *transition matrix* of this time-homogeneous Markov chain is

$$P_{i,j} := \mathbb{P}(B_k = j | B_{k-1} = i) = \begin{cases} \frac{1}{\text{deg}(i)} & : i \sim j \\ 0 & : \text{otherwise.} \end{cases}$$

Denote by $(L_k)_{k \geq 0}$ the random walk lifted to the node attributes in \mathbb{R}^d , that is

$$L_k := f(B_k). \quad (6)$$

Recall that the *normalized graph Laplacian* for random walks is defined as $\mathcal{L} = I - D^{-1}A$, where D is diagonal degree matrix; in particular, the entry-wise definition is

$$\mathcal{L}_{i,j} := \begin{cases} -\frac{1}{\text{deg}(i)} & : i \sim j \\ 1 & : i = j \\ 0 & : \text{otherwise.} \end{cases}$$

The discrete graph diffusion equation for $U_k \in \mathbb{R}^{n \times d}$ is given by

$$U_k - U_{k-1} = -\mathcal{L}U_{k-1}, \quad U_0^{(i)} = f(i) \quad (7)$$

³The labels given by the labelled graph are called *attributes*, while the computed updates are called *features*.

where the initial condition $U_0 \in \mathbb{R}^{n \times d}$ is specified by the node attributes.⁴ The probabilistic interpretation of the solution to this diffusion equation is classical and given as

$$U_k = (\mathbb{E}[L_k | B_0 = i])_{i=1}^n = P^k U_0. \quad (8)$$

This allows us to compute the solution u_k using the transition matrix $P = I - \mathcal{L}$.

Random Walks on Algebras. We now incorporate the history of a random walker by considering the quantity

$$\mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i] = \mathbb{E}[\varphi(\delta_0 \mathbf{L}) \cdots \varphi(\delta_k \mathbf{L}) | B_0 = i] \quad (9)$$

where $\mathbf{L}_k = (L_0, \dots, L_k)$. Since $\tilde{\varphi}$ captures the whole history of the random walk \mathbf{L}_k over node attributes, we expect this expectation to provide a much richer summary of the neighborhood of node i than $\mathbb{E}[L_k | B_0 = i]$. The price is however, the computational complexity, since (9) is H -valued. We first show, that analogous to (7), the quantity (9) satisfies a diffusion equation that can be computed with linear algebra. To do so, we develop a graph analogue of the hypo-elliptic Laplacian and replace the scalar entries of the matrices with entries from the algebra H .

Matrix Rings over Algebras. We first revisit the adjacency matrix $A \in \mathbb{R}^{n \times n}$ and replace it by the *tensor adjacency matrix* $\tilde{A} = (\tilde{A})_{i,j} \in H^{n \times n}$, that is \tilde{A} is a matrix but instead of scalar entries its entries are elements in the algebra H . The matrix A has an entry at i, j if nodes i and j are connected; \tilde{A} replaces the i, j entry with an element of H that tells us how the node attributes of i and j differ,

$$\tilde{A}_{i,j} := \begin{cases} \varphi(f(j) - f(i)) & : i \sim j \\ 0 & : \text{otherwise.} \end{cases} \quad (10)$$

Matrix multiplication works for elements of $H^{n \times n}$ by replacing scalar multiplication by multiplication in H , that is $(\tilde{B} \cdot \tilde{C})_{i,j} = \sum_{k=1}^n \tilde{B}_{i,k} \cdot \tilde{C}_{k,j}$ for $\tilde{B}, \tilde{C} \in H^{n \times n}$ and $\tilde{B}_{i,k} \cdot \tilde{C}_{k,j}$ denotes multiplication in H as in Equation (2). For the classical adjacency matrix A , the k -th power counts the number of length k walks in the graph, so that $(A^k)_{i,j}$ is the number of walks of length k from node i to node j .

We can take powers of \tilde{A} in the same way as in the classical case, where

$$(\tilde{A}^k)_{i,j} = \sum_{\mathbf{x}} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x})$$

where the sum is taken over all length k walks $\mathbf{x} = (f(i), \dots, f(j))$ from node i to node j (full details are provided in Appendix D). Since $\tilde{\varphi}(\mathbf{x})$ characterizes each walk \mathbf{x} , the entry $\tilde{A}_{i,j}^k$ can be interpreted as a summary of all walks which connect nodes i and j .

Hypo-elliptic Graph Diffusion. Similar to the tensor adjacency matrix, we define the *hypo-elliptic graph Laplacian* as the $n \times n$ matrix

$$\tilde{\mathcal{L}} = I - D^{-1} \tilde{A} \in H^{n \times n},$$

where D is the degree matrix embedded into $H^{n \times n}$ at tensor degree 0. The entry-wise definition is

$$\tilde{\mathcal{L}}_{i,j} := \begin{cases} \frac{-\varphi(f(j)-f(i))}{\deg(i)} & : i \sim j \\ 1 & : i = j \\ 0 & : \text{otherwise.} \end{cases} \quad (11)$$

We can now formulate the *hypo-elliptic graph diffusion equation* for $\mathbf{v}_k \in H^n$ as

$$\mathbf{v}_k - \mathbf{v}_{k-1} = -\tilde{\mathcal{L}} \mathbf{v}_{k-1}, \quad \mathbf{v}_0^{(i)} = \varphi(f(i)). \quad (12)$$

Analogous to the classic graph diffusion (8), the hypo-elliptic graph diffusion (12) has a probabilistic interpretation in terms of \mathbf{L} as shown in Theorem 1 (the proof is given in Appendix D).

Theorem 1. *Let $k \in \mathbb{N}$, $\mathbf{L}_k = (L_0, \dots, L_k)$ be the lifted random walk from (6), and $\tilde{P} = I - \tilde{\mathcal{L}}$ be the tensor adjacency matrix. The solution to the hypo-elliptic graph diffusion equation (12) is*

$$\mathbf{v}_k = (\mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_0 = i])_{i=1}^n = \tilde{P}^k \mathbf{1}_H.$$

Furthermore, if $F \in H^{n \times n}$ is the diagonal matrix with $F_{i,i} = \varphi(f(i))$, then

$$F \mathbf{v}_k = (\mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i])_{i=1}^n.$$

⁴The attributes over all nodes are given by an $n \times d$ matrix; in particular $U_k^{(i)}$ is the i^{th} row of the matrix.

In the classical diffusion equation, U_k captures the concentration of the random walkers after k time steps over the nodes. In the hypo-elliptic diffusion equation, \mathbf{v}_k captures summaries of random walk histories after k time steps over the nodes since $\tilde{\varphi}(\mathbf{L}_k)$ summarizes the whole trajectory $\mathbf{L}_k = (L_0, \dots, L_k)$ and not only the endpoint L_k .

Node Features and Graph Features. Theorem 1 can then be used to compute features $\Phi(i) \in H$ for individual nodes as well as a feature $\Psi(\mathcal{G})$ for the entire graph. The former is given by i -th component $\mathbf{v}_k^{(i)}$ of the solution $\mathbf{v}_k = (\mathbf{v}_k^{(i)})_{i=1, \dots, n} \in H^n$ of Equation (12),

$$\Phi(i) := \mathbf{v}_k^{(i)} = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i] = (F\tilde{P}^k \mathbf{v}_0)^{(i)} \in H,$$

since the random walk B chooses the starting node $B_0 = i$ uniformly at random. The latter can be computed by mean pooling the node features, which also has a probabilistic interpretation as

$$\Psi(\mathcal{G}) := \frac{1}{n} \sum_{i=1}^n \mathbf{v}_k^{(i)} = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k)] = n^{-1} (\mathbb{1}_H^T F \tilde{P}^k \mathbf{v}_0) \in H, \quad (13)$$

where $\mathbb{1}_H^T := (1_H, \dots, 1_H) \in H^n$ is the all-ones vector in H and 1_H denotes the unit in H .

Characterizing Graphs with Random Walks. The graph and node features obtained through the hypo-elliptic diffusion equation are highly descriptive: they characterize the entire history of the random walk process if one also includes the time parametrization, as described in Appendix C.

Theorem 2. *Suppose Ψ is the graph feature map from Equation (13) induced by the tensor exponential algebra lifting including time parametrization. Let \mathcal{G} and \mathcal{G}' be two labelled graphs, and $\mathbf{L}_k = (L_0, \dots, L_k)$ and $\mathbf{L}'_k = (L'_0, \dots, L'_k)$ be the k -step lifted random walk as defined in Equation (6). Then, $\Psi(\mathcal{G}) = \Psi(\mathcal{G}')$ if and only if the distributions of \mathbf{L}_k and \mathbf{L}'_k are equal.*

It is instructive to contrast this result with the classical diffusion case; the latter only uses the marginal distribution of L_k to capture the graph structure, which at least intuitively has much less expressive power. Indeed, in Appendix E, we show that for elementary graphs, this already leads to big differences in expressive power. Further, an analogous result holds for the node features, and we prove both results in Appendix E. While we use the tensor exponential in this article, many other choices of $\tilde{\varphi}$ are possible and result in graph and node features with such properties: under mild conditions, if the algebra lifting $\varphi : \mathbb{R}^d \rightarrow H$ characterizes measures on \mathbb{R}^d , the resulting node feature map Φ characterizes the random walk, see [62], which in turn implies the above results. Possible variations are discussed in Appendix G.

General (Hypo-elliptic) Diffusions and Attention. One can consider more general diffusion operators, such as the normalized Laplacian \mathcal{K} of a weighted graph. We define its lifted operator $\tilde{\mathcal{K}} \in H^{n \times n}$ analogous to Equation (11), resulting in a generalization of Theorem 1 with $\tilde{\mathcal{K}}$ replacing $\tilde{\mathcal{L}}$. In the flavour of convolutional GNNs [8], we consider a weighted adjacency matrix $A \in \mathbb{R}^{n \times n}$

$$A_{i,j} = \begin{cases} c_{i,j} & : i \sim j \\ 0 & : \text{otherwise,} \end{cases}$$

for $c_{i,j} > 0$. The corresponding normalized Laplacian \mathcal{K} is given by $\mathcal{K} = I - D^{-1}A$, where D is a diagonal matrix with $D_{i,i} = \sum_{j \in \mathcal{N}(i)} c_{i,j}$. A common way to learn the coefficients is by introducing parameter sharing across graphs by modelling them as $c_{i,j} = \exp(a(f(i), f(j)))$ using a local attention mechanism, $a : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ [65]. In our implementation, we use additive attention [4] given by $a(f(i), f(j)) = \text{LeakyReLU}_{0.2}(W_s f(i) + W_t f(j))$, where $W_s, W_t \in \mathbb{R}^{1 \times d}$ are linear transformations for the source and target nodes, but different attention mechanisms can also be used; e.g. scaled dot-product attention [64]. Then, the corresponding transition matrix $P = D^{-1}A$ is defined as $P_{ij} = \text{softmax}_{k \in \mathcal{N}(i)}(a(f(i), f(k)))_j$. The lifted transition matrix is defined as

$$\tilde{P} = \begin{cases} P_{i,j} \varphi(f(j) - f(i)) & : i \sim j \\ 0 & : \text{otherwise.} \end{cases}$$

The statements of Theorem 1 immediately generalize to this variation by replacing the expectation with respect to a non-uniform random walk. Hence, in this case the use of attention can be interpreted as learning the transition probabilities of a random walk on the graph.

4 Efficient Algorithms for Deep Learning

The previous sections show that the node feature $\Phi(i)$ provides a structured description of the neighborhood of node i and it is instructive to think of a linear functional $\langle \ell, \Phi(i) \rangle$ as answering a specific question about the node neighbourhood, see Appendix E for examples. The naive computation of $\langle \ell, \Phi(i) \rangle$ by first computing $\Phi(i)$ and taking the inner product is too expensive, especially when $\ell = (\ell_0, \dots, \ell_M, 0, \dots) \in H$ for large M . To address this we revisit two observations from [62]: first, for a rank-1 functional $\ell \in H$, the computation of $\langle \ell, \Phi(i) \rangle$ is computationally cheap. Second, restriction to small M limits the expressive power but can be counteracted by composition: any choice of d different functionals $\ell^1, \dots, \ell^d \in H$ gives a label update $f(i) \mapsto (\langle \ell^j, \Phi(i) \rangle)_{j=1, \dots, d} \in \mathbb{R}^d$ for the graph. Repeating such a label update a few times with low-degree M and rank-1 functionals turns out to be as powerful as computing one update for general functionals with arbitrary high M . The first observation should not be too surprising given the popularity of low rank approximations; the second observation is reminiscent to constructing a high-degree polynomial by composing low-degree polynomials⁵ or the width-vs-depth phenomenon in neural nets and we give more details below.

Computing Rank-1 Functionals. First, we focus on a *rank-1* linear functional $\ell \in H$ given as

$$\ell = (\ell_m)_{m \geq 0} \text{ with } \ell_m = u_{M-m+1} \otimes \dots \otimes u_M \text{ and } \ell_m = 0 \text{ for } m > M, \quad (14)$$

where $u_m \in \mathbb{R}^d$ for $m = 1, \dots, M$ for a fixed $M \geq 1$. Theorem 3 shows that for such ℓ , the computation of $\langle \ell, \hat{\Phi}(i) \rangle$, where $\hat{\Phi}(i)$ is the node feature without the basepoint, can be done (a) efficiently by factoring this low-rank structure into the recursive computation, and (b) simultaneously for all nodes $i \in \mathcal{V}$ in parallel. This can then be used to compute rank- R functionals for $R > 1$, and for $\langle \ell, \Phi(i) \rangle$; see Appendix F, where we also provide a pseudocode implementation.

Theorem 3. Let ℓ be as in (14) and define $f_{k,m} \in \mathbb{R}^n$ for $m = 1, \dots, M$ as

$$f_{1,m} := \frac{1}{m!} (P \odot C^{u_{M-m+1}} \odot \dots \odot C^{u_M}) \cdot \mathbb{1},$$

where $\mathbb{1}^T := (1, \dots, 1) \in \mathbb{R}^n$ is the all-ones vector; and for $2 \leq k$ and $1 \leq m \leq M$ recursively as

$$f_{k,m} := P \cdot f_{k-1,m} + \sum_{r=1}^m \frac{1}{r!} (P \odot C^{u_{M-m+1}} \odot \dots \odot C^{u_{M-m+r}}) \cdot f_{k-1,m-r}, \quad (15)$$

where the matrix $C^u = (C_{i,j}^u) \in \mathbb{R}^{n \times n}$ is defined as

$$C_{i,j}^u := \begin{cases} \langle u, f(j) - f(i) \rangle & : i \sim j, \\ 0 & : \text{otherwise.} \end{cases}$$

Here \odot denotes element-wise⁶ multiplication, while \cdot denotes matrix multiplication. Then, it holds for $i \in \mathcal{V}$, random walk length $k \in \mathbb{Z}_+$, and tensor degree $m = 1, \dots, M$, that

$$f_{k,m}(i) = \langle \ell_m, \hat{\Phi}_k(i) \rangle,$$

where $\hat{\Phi}_k(i) = \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \dots \varphi(\delta_k \mathbf{L}_k) \mid B_0 = i]$.

Overall, Eq. (15) computes $f_{k,m}(i)$ for all $i \in \mathcal{V}$, $k = 1, \dots, K$, $m = 1, \dots, M$ in $O(K \cdot M^2 \cdot N_E + M \cdot N_E \cdot d)$ operations, where $N_E \in \mathbb{N}$ denotes the number of edges; see App. F. In particular, one does not need to compute $\Phi(i) \in H$ directly or store large tensors.

Graph Labelling Layers. Fixing d rank 1-functionals $\ell^1, \dots, \ell^d \in H$ induces a label update $f(i) \mapsto (\langle \ell^i, \Phi(i) \rangle)_{i=1, \dots, d} \in \mathbb{R}^d$. Theorem 3 allows us to compute this update in parallel for all nodes in \mathcal{V} . Such a label update is similar to hidden layer in a NN and we can stack such updates, see Figure 3 in App. H.1. As in NN, the d functionals in each "graph labelling layer" are optimized by gradient descent. Finally, note that a rank R functional is the sum of R rank-1 functionals so we can immediately carry out the same construction with rank- R functionals by adding a mixing layer.

⁵For example, $1 + x + x^2$ composed with $1 + 2x^2$ yields the degree 4 polynomial $1 + (1 + 2x^2) + (1 + 2x^2)^2$.

⁶For example $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \odot \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 5 & 12 \\ 21 & 32 \end{bmatrix}$.

To sum up, a graph labelling layer is determined by the random walk length k , the maximal tensor degree M , maximal tensor rank R and the functionals are then found by optimization.

Using a single layer of low-rank functionals limits the expressiveness but stacking layers allows in practice to approximate general, high-degree M functionals. Some theoretical results can be found in [62]; however, here we simply appeal to the analogy with NN where stacking simple transformations provides a flexible functional class with good inductive bias.

5 Experiments

We implemented the above approach and call the resulting model **Graph2Tens Networks** since it represents the neighbourhood of a node as a sequence of tensors, which is further pushed through a low-tensor-rank constrained linear mapping, similarly to how neural networks linearly transform their inputs pre-activation. A conceptual difference is that in our case the non-linearity is applied first and the projection secondly, albeit the computation is coupled between these steps. We provide further experiments and ablation studies of our models in Appendix H.

Experimental Setup. The aim of our main experiment is to test the following key properties of our model: (1) ability to capture long-range interactions between nodes in a graph, (2) robustness to pooling operations, hence making it less susceptible to the “over-squashing” phenomenon [2]. We do this by following the experiments in [70]. In particular, we show that our model is competitive with previous approaches for retaining long-range context in graph-level learning tasks but without computing all pairwise interactions between nodes, thus keeping the influence distribution localized [73]. We further give a detailed ablation study to show the robustness of our model to various architectural choices in Appendix H.2. As a second experiment, we follow the previous applications of diffusion approaches to graphs that have mostly considered inductive learning tasks, e.g. on the citation datasets [13, 60, 11]. Our experimentation on these datasets are available in Appendix H.3, where the model performs on par with short-range GNN models, but does not seem to benefit from added long-range information a-priori. However, when labels are dropped in a k -hop sanitized way as in [50], the performance decrease is less pronounced.

Datasets. We use two biological graph classification datasets (NCI1 and NCI109), that contain around ~ 4000 biochemical compounds represented as graphs with ~ 30 nodes on average [67, 1]. The task is to predict whether a compound contains anti-lung-cancer activity. The dataset is split in a ratio of 80% – 10% – 10% for training, validation and testing. Previous work [2] has found that GNNs that only summarize local structural information can be greatly outperformed by models that are able to account for global contextual relationships through the use of *fully-adjacent* layers. This was further improved on by [70], where a local neighbourhood encoder consisting of a GNN stack was upgraded with a Transformer submodule [64] for learning global interactions.

Model Details. We build a GNN architecture primarily motivated by the GraphTrans (small) model from [70], and only fine-tune the pre- and postprocessing layers(s), random walk length, functional degree and optimization settings. In detail, a preprocessing MLP layer with 128 hidden units is followed by a stack of 4 G2TN layers each with RW length-5, max rank-128, max tensor degree-2, all equipped with JK-connections [73] and a max aggregator. Afterwards, the node features are combined into a graph-level representation using gated attention pooling [41]. The pooled features are transformed using a final MLP layer with 256 hidden units, and then fed into a softmax classification layer. The pre- and postprocessing MLP layers employ skip-connections [30]. Both MLP and G2TN layers are followed by layer normalization [3], where GTN layers normalize their rank-1 functionals independently across different tensor degrees, which corresponds to a particular realization of group normalization [69]. We randomly drop 10% of the features for all hidden layers during training [58]. The attentional variant, G2T(A)N also randomly drops 10% of its edges and uses 8 attention heads [65]. Training is performed by minimizing the categorical cross-entropy loss with an ℓ_2 regularization penalty of 10^{-4} . For optimization, Adam [32] is used with a batch size of 128 and an initial learning rate of 10^{-3} that is decayed via a cosine annealing schedule [42] over 200 epochs. Further intuition about the model and architectural choices are available in Appendix H.1.

Baselines. We compare against (1) the baseline models reported in [70], (2) variations of GraphTrans, (3) other recently proposed hierarchical approaches for long-range graph tasks [50]. Groups

Table 1: Comparison of classification accuracies on NCI biological datasets, where we report mean and standard deviation over 10 random seeds for our models.

Model	GNN Type	GNN Count	NCI1 (%)	NCI109 (%)
Set2Set [40, 66]	GCN	3	68.6 ± 1.9	69.8 ± 1.2
SortPool [40, 75]	GCN	3	73.8 ± 1.0	74.0 ± 1.2
SAGPool _h [40]	GCN	3	67.5 ± 1.1	67.9 ± 1.4
SAGPool _g [40]	GCN	3	74.2 ± 1.2	74.1 ± 0.8
GIN [21, 72]	GIN	8	80.0 ± 1.4	-
GCN + VN [74, 24]	GCN	2	71.5	-
HGNet-EdgePool [74, 54]	GCN+RGCN	3 + 2	77.1	-
HGNet-Louvain [74, 54]	GCN+RGCN	3 + 2	75.1	-
GIN + FA [2, 72]	GIN	8	81.5 ± 1.2	-
GraphTrans (small) [70, 64]	GCN	3	81.3 ± 1.9	79.2 ± 2.2
GraphTrans (large) [70, 64]	GCN	4	82.6 ± 1.2	82.3 ± 2.6
G2T(A)N (ours)	G2T(A)N	4	81.9 ± 1.2	78.0 ± 2.3
G2TN (ours)	G2TN	4	80.7 ± 2.5	78.9 ± 2.5

of models in Table 1 are separated by dashed lines if they were reported in separate papers, and the first citation after the name is where the result first appeared. The number of GNN layers in HGNet are not discussed by [50], and we report it as implied by their code. We organize the models into three groups divided by solid lines: (a) baselines that only apply neighbourhood aggregations, and hierarchical or global pooling schemes, (b) baselines that first employ a local neighbourhood encoder, and afterwards fully densify the graph in one way or another so that all nodes interact with each other *directly*, (c) our models that we emphasize thematically belong to (a).

Results. In Table 1, we report the mean and standard deviation of classification accuracy computed over 10 different seeds. Overall, both our models improve over all baselines in group (a) on both datasets, maximally by 1.9% on NCI1 and by 4.8% on NCI109. In group (b), G2T(A)N is solely outperformed by GraphTrans (large) on NCI1 by only 0.7%. Interestingly, the attention-free variation, G2TN, performs better on NCI109, where it performs very slightly worse than GraphTrans (small).

Discussion. The previous experiments demonstrate that our approach performs very favourably on long-range reasoning tasks compared to GNN-based alternatives without global pairwise node interactions. Several of the works we compare against have focused on extending GNNs to larger neighbourhoods by specifically designed graph coarsening and pooling operations, and we emphasize two important points: (1) our approach can efficiently capture large neighbourhoods without any need for coarsening, (2) it already performs well with simple mean-pooling as justified by Theorem 2 and experimentally supported by the ablation studies in Appendix H.2. Although the Transformer-based GraphTrans slightly outperforms our model potentially due to its ability to learn global interactions, it is not entirely clear how much of the global graph structure it is able to infer from interactions of short-range neighbourhood summaries. Finally, Transformer models can be bottlenecked by their quadratic complexity in nodes, while our approach only scales with edges, and hence, it can be more favourable for large sparse graphs in terms of computations.

6 Conclusion

Inspired by classical results from analysis [23], we introduce the hypo-elliptic graph Laplacian. This yields a diffusion equation and also generalizes its classical probabilistic interpretation via random walks but now taking history into account. In addition to several attractive theoretical guarantees, we provide scalable algorithms. Our experiments show that this can lead to largely improved baselines for long-range reasoning tasks. A promising future research theme is to develop improvements for the classical Laplacian in this hypo-elliptic context; including lazy random walks [71]; nonlinear diffusions [12]; and source/sink terms [60]. Another theme could be to extend the geometric study [61] of over-squashing to this hypo-elliptic point of view which is naturally tied to sub-Riemannian geometry [59]. A limitation in our theoretical results is that for the iterations of low-rank approximations only partial results exist and expanding this is an interesting (algebra-heavy) topic.

Acknowledgments and Disclosure of Funding

Csaba Toth was supported by a Mathematical Institute Award from the University of Oxford. Celia Hacker and Darrick Lee were supported by NCCR-Synapsy Phase-3 SNSF grant number 51NF40-185897. Darrick Lee and Harald Oberhauser were supported by the Hong Kong Innovation and Technology Commission (InnoHK Project CIMDA). Harald Oberhauser was also supported by the DataSig Program [EP/S026347/1] and the Oxford-Man Institute.

References

- [1] Pubchem. <https://pubchem.ncbi.nlm.nih.gov/>.
- [2] Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical implications. In *International Conference on Learning Representations*, 2021.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [4] Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- [5] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. *Advances in neural information processing systems*, 29, 2016.
- [6] Patric Bonnier, Chong Liu, and Harald Oberhauser. Adapted topologies and higher rank signatures. *arXiv preprint arXiv:2005.08897*, 2020.
- [7] Karsten Borgwardt, Elisabetta Ghisu, Felipe Llinares-López, Leslie O’Bray, and Bastian Rieck. Graph kernels: State-of-the-art and future challenges. *arXiv preprint arXiv:2011.03854*, 2020.
- [8] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges. *arXiv:2104.13478 [cs, stat]*, May 2021.
- [9] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- [10] J. Douglas Carroll and Jih-Jie Chang. Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition. *Psychometrika*, 35(3):283–319, September 1970.
- [11] Benjamin Paul Chamberlain, James Rowbottom, Davide Eynard, Francesco Di Giovanni, Xiaowen Dong, and Michael M. Bronstein. Beltrami flow and neural diffusion on graphs. *CoRR*, abs/2110.09443, 2021.
- [12] Benjamin Paul Chamberlain, James Rowbottom, Maria I. Gorinova, Stefan D. Webb, Emanuele Rossi, and Michael M. Bronstein. GRAND: Graph Neural Diffusion. In *The Symbiosis of Deep Learning and Differential Equations*, September 2021.
- [13] Benjamin Paul Chamberlain, James R. Rowbottom, Maria I. Gorinova, Stefan Webb, Emanuele Rossi, and Michael M. Bronstein. Grand: Graph neural diffusion. In *ICML*, 2021.
- [14] Dexiong Chen, Laurent Jacob, and Julien Mairal. Convolutional kernel networks for graph-structured data. In *International Conference on Machine Learning*, pages 1576–1586. PMLR, 2020.
- [15] Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. *Advances in neural information processing systems*, 32, 2019.
- [16] Ilya Chevyrev and Harald Oberhauser. Signature moments to characterize laws of stochastic processes. *arXiv:1810.10971 [math, stat]*, 2018.

- [17] Thomas Cochrane, Peter Foster, Varun Chhabra, Maud Lemercier, Terry Lyons, and Cristopher Salvi. Sk-tree: a systematic malware detection algorithm on streaming trees via the signature kernel. In *2021 IEEE International Conference on Cyber Security and Resilience (CSR)*, pages 35–40. IEEE, 2021.
- [18] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- [19] Persi Diaconis. Group representations in probability and statistics. *Lecture notes-monograph series*, 11:i–192, 1988.
- [20] Ahmed A. A. Elhag, Gabriele Corso, Hannes Stärk, and Michael M. Bronstein. Graph anisotropic diffusion, 2022.
- [21] Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph neural networks for graph classification. In *International Conference on Learning Representations*, 2019.
- [22] Michel Fliess. Fonctionnelles causales non linéaires et indéterminées non commutatives. *Bulletin de la société mathématique de France*, 109:3–40, 1981.
- [23] Bernard Gaveau. Principe de moindre action, propagation de la chaleur et estimatees sous elliptiques sur certains groupes nilpotents. *Acta Mathematica*, 139(none):95–153, January 1977.
- [24] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.
- [25] Daniele Grattarola and Cesare Alippi. Graph neural networks in tensorflow and keras with spektral [application notes]. *IEEE Computational Intelligence Magazine*, 16(1):99–106, 2021.
- [26] Alexander Grigoryan. *Heat kernel and analysis on manifolds*, volume 47. American Mathematical Soc., 2009.
- [27] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.
- [28] Ben Hambly and Terry Lyons. Uniqueness for the signature of a path of bounded variation and the reduced path group. *Ann. of Math.*, 171(1):109–167, 2010.
- [29] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [31] Lars Hörmander. Hypoelliptic second order differential equations. *Acta Mathematica*, 119:147–171, 1967.
- [32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [33] Thomas Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations*, 2017.
- [34] Franz J. Kiraly and Harald Oberhauser. Kernels for sequentially ordered data. *J. Mach. Learn. Res.*, 20(31):1–45, 2019.
- [35] Johannes Klicpera, Stefan Weissenberger, and Stephan Günnemann. Diffusion improves graph learning. *ArXiv*, abs/1911.05485, 2019.
- [36] Risi Kondor. Diffusion kernels on graphs and other discrete structures. In *ICML 2002*, 2002.

- [37] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- [38] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [39] Darrick Lee and Robert Ghrist. Path Signatures on Lie Groups. *arXiv:2007.06633 [cs, math, stat]*, 2020.
- [40] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *International conference on machine learning*, pages 3734–3743. PMLR, 2019.
- [41] Yujia Li, Richard Zemel, Marc Brockschmidt, and Daniel Tarlow. Gated graph sequence neural networks. In *Proceedings of ICLR’16*, 2016.
- [42] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017.
- [43] Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. *Advances in neural information processing systems*, 32, 2019.
- [44] Ron Milo, Shai S. Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri B. Chklovskii, and Uri Alon. Network motifs: simple building blocks of complex networks. *Science*, 298 5594:824–7, 2002.
- [45] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5115–5124, 2017.
- [46] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):4602–4609, July 2019.
- [47] Galileo Namata, Ben London, Lise Getoor, Bert Huang, and U Edu. Query-driven active surveying for collective classification. In *10th International Workshop on Mining and Learning with Graphs*, volume 8, page 1, 2012.
- [48] Giannis Nikolentzos and Michalis Vazirgiannis. Random walk graph neural networks. *Advances in Neural Information Processing Systems*, 33:16211–16222, 2020.
- [49] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710, 2014.
- [50] Ladislav Rampásek and Guy Wolf. Hierarchical graph neural nets can capture long-range interactions. In *2021 IEEE 31st International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE, 2021.
- [51] Linda Preiss Rothschild and Elias M. Stein. Hypoelliptic differential operators and nilpotent groups. *Acta Mathematica*, 137:247–320, 1976.
- [52] Alaa Saade, Florent Krzakala, and Lenka Zdeborová. Spectral clustering of graphs with the bethe hessian. In *NIPS*, 2014.
- [53] 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, 2008.
- [54] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *European semantic web conference*, pages 593–607. Springer, 2018.
- [55] Alice C Schwarze and Mason A Porter. Motifs for processes on networks. *SIAM Journal on Applied Dynamical Systems*, 20(4):2516–2557, 2021.

- [56] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- [57] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- [58] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [59] Robert S Strichartz. Sub-riemannian geometry. *Journal of Differential Geometry*, 24(2):221–263, 1986.
- [60] Matthew Thorpe, Tan Minh Nguyen, Hedi Xia, Thomas Strohmer, Andrea Bertozzi, Stanley Osher, and Bao Wang. GRAND++: Graph neural diffusion with a source term. In *International Conference on Learning Representations*, 2022.
- [61] Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature, 2021.
- [62] Csaba Toth, Patric Bonnier, and Harald Oberhauser. Seq2Tens: An Efficient Representation of Sequences by Low-Rank Tensor Projections. In *International Conference on Learning Representations*, 2021.
- [63] N Th Varopoulos. Analysis and geometry on groups. *Cambridge Tracts in Math.*, 100, 1992.
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [65] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [66] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets. In *ICLR*, 2016.
- [67] Nikil Wale, Ian A. Watson, and George Karypis. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and Information Systems*, 14(3):347–375, March 2008.
- [68] Wolfgang Woess. *Random Walks on Infinite Graphs and Groups*. Cambridge Tracts in Mathematics. Cambridge University Press, Cambridge, 2000.
- [69] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.
- [70] Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E Gonzalez, and Ion Stoica. Representing long-range context for graph neural networks with global attention. *Advances in Neural Information Processing Systems*, 34:13266–13279, 2021.
- [71] Louis-Pascal A. C. Xhonneux, Meng Qu, and Jian Tang. Continuous graph neural networks. In *Proceedings of the 37th International Conference on Machine Learning, ICML’20*, pages 10432–10441. JMLR.org, July 2020.
- [72] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [73] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *International Conference on Machine Learning*, pages 5453–5462. PMLR, 2018.

- [74] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, pages 4805–4815, Red Hook, NY, USA, December 2018. Curran Associates Inc.
- [75] Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep learning architecture for graph classification. In *Thirty-second AAAI conference on artificial intelligence*, 2018.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) While our approach is less susceptible to the “over-squashing” phenomenon and empirically performs well on learning long-range interactions, we find that our methods are less beneficial for inductive learning tasks. This is stated explicitly in Section 5, with experiments shown in Appendix H.3.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[No\]](#) This article introduce a method which can improve the performance on graph neural networks. We don’t think there are any foreseeable negative societal impacts beyond generic points that apply to any graph learning method.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) All theorem statements contain the full set of assumptions.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) Proofs to all of our theoretical results are contained in the appendix, with specific references to sections within the main text.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) , see Appendix H
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) This is detailed in Section 5, with further discussion in Appendix H.2.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) All experiments are run using multiple seeds, and the standard deviation over these seeds is provided.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Appendix H.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
 - (b) Did you mention the license of the assets? [\[Yes\]](#) We include it in our code.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#) We include code as part of supplementary material.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Broader Impacts

This article introduces a method which can improve the performance on graph neural networks. We do not believe there are any foreseeable negative societal impacts beyond generic points that apply to any graph learning method.

Appendix Outline

Section A summarizes the notation and objects used in this paper. Section B gives general background on tensor and the algebra H . Section C discusses the feature map $\tilde{\varphi}$ for sequences and possible choices for the lift φ from labels to the algebra H . Section D contains the proofs of our main theorems and some variations of hypoelliptic diffusion. Section E provides the formal statement of Theorem 2 and the extension to pooled features. Section F gives background and details of the low-rank algorithm. Section G discusses variations of the sequence feature map which lead to different node and graph features. Section H includes further experiments and discussion on the empirical results.

A Notation

Symbol	Meaning
Fixed Parameters and Indices	
d	dimension of node attributes
k	length of random walk
n	number of nodes in graph
Sequence Features	
\mathbb{R}^d	finite dimensional vector space for node attributes
$\text{Seq}(\mathbb{R}^d)$	sequences $\mathbf{x} = (x_0, \dots, x_k)$ of arbitrary length k in \mathbb{R}^d
$\delta_k \mathbf{x}$	increments of a sequence where $\delta_0 \mathbf{x} := x_0$ and $\delta_k \mathbf{x} := x_k - x_{k-1}$ for $k > 0$
H	tensor algebra of \mathbb{R}^d (see Appendix B)
φ	algebra lifting $\varphi : \mathbb{R}^d \rightarrow H$
\exp_{\otimes}	tensor exponential $\exp_{\otimes} : \mathbb{R}^d \rightarrow H$ (main example of algebra lifting)
$\tilde{\varphi}$	sequence feature map $\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H$, where $\tilde{\varphi}(\mathbf{x}) = \varphi(\delta_0 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x})$
Graphs, Adjacency and Laplacian Matrices	
\mathcal{G}	$\mathcal{G} = (\mathcal{V}, \mathcal{E}, f)$ graph with vertex set, edge set, and node attributes $f : \mathcal{V} \rightarrow \mathbb{R}^d$
$(B_k)_{k \geq 0}$	simple random walk on graph (valued in \mathcal{V})
$(L_k)_{k \geq 0}$	lifted random walk over \mathbb{R}^d , $L_k := f(B_k)$
\mathbf{L}_k	lifted length k random walk $\mathbf{L}_k = (L_0, \dots, L_k)$
Φ	$\Phi(i) = \mathbb{E}[\tilde{\varphi}(L_1, \dots, L_n) L_0 = i] \in H$ feature map for node $i \in \mathcal{V}$
Ψ	$\Psi(\mathcal{G}) = \mathbb{E}[\tilde{\varphi}(L_1, \dots, L_n)] \in H$ feature map for graphs
A	standard adjacency matrix
\tilde{A}	tensor adjacency matrix
P	standard transition matrix
\tilde{P}	tensor transition matrix
\mathcal{L}	normalized graph Laplacian
$\tilde{\mathcal{L}}$	normalized hypo-elliptic graph Laplacian

Notation Conventions

- Bold symbols are used for tensors and sequences; vectors are unbolded, such as $x \in \mathbb{R}^d$
- Coordinates for vectors are denoted using superscripts: $x = (x^{(1)}, \dots, x^{(d)}) \in \mathbb{R}^d$.
- Tensors are denoted by $\mathbf{v} = (\mathbf{v}_0, \mathbf{v}_1, \mathbf{v}_2, \dots) \in H$, where $\mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}$.
- Sequences are denoted by $\mathbf{x} = (x_0, \dots, x_k) \in \text{Seq}(\mathbb{R}^d)$.

B Tensors and the Algebra H

In this section, we provide a brief overview of tensors on a finite-dimensional vector space \mathbb{R}^d , along with the resulting algebra H .

Tensors on \mathbb{R}^d . While tensor products between vector spaces are defined more generally, our main focus is on defining the tensor powers of \mathbb{R}^d , denoted by $(\mathbb{R}^d)^{\otimes m}$ for some $m \in \mathbb{N}$. The tensor power $(\mathbb{R}^d)^{\otimes m}$ is also a vector space. Given a basis e^1, \dots, e^d of \mathbb{R}^d , we can define a basis of $(\mathbb{R}^d)^{\otimes m}$ as the collection of all

$$e^I := e^{i_1} \otimes \dots \otimes e^{i_m}$$

over all *multi-indices* (i_1, \dots, i_m) , where each $i_j \in [d]$. Thus, any element $\mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}$ can be represented as

$$\mathbf{v}_m = \sum_{i_1, \dots, i_m=1}^d \mathbf{v}_m^{(i_1, \dots, i_m)} e^{(i_1, \dots, i_m)},$$

where the $\mathbf{v}_m^{(i_1, \dots, i_m)} \in \mathbb{R}$ specifies the coordinates of \mathbf{v}_m . In particular, $(\mathbb{R}^d)^{\otimes 1}$ is simply \mathbb{R}^d itself; $(\mathbb{R}^d)^{\otimes 2}$ can be viewed as the space of $d \times d$ matrices; $(\mathbb{R}^d)^{\otimes 3}$ is the space of $d \times d \times d$ arrays, etc.

Given two vectors $x, y \in \mathbb{R}^d$ which we represent coordinate-wise as

$$x = (x^{(1)}, \dots, x^{(d)}), \quad y = (y^{(1)}, \dots, y^{(d)}),$$

the tensor product $x \otimes y \in (\mathbb{R}^d)^{\otimes 2}$ is given by

$$(x \otimes y)^{(i,j)} := x^{(i)} y^{(j)}.$$

More generally, if $\mathbf{u}_m \in (\mathbb{R}^d)^{\otimes m}$ and $\mathbf{v}_n \in (\mathbb{R}^d)^{\otimes n}$, the tensor product $\mathbf{u}_m \otimes \mathbf{v}_n \in (\mathbb{R}^d)^{\otimes (m+n)}$ is defined coordinate-wise by

$$(\mathbf{u}_m \otimes \mathbf{v}_n)^{(i_1, \dots, i_{m+n})} := \mathbf{u}_m^{(i_1, \dots, i_m)} \mathbf{v}_n^{(i_{m+1}, \dots, i_{m+n})}.$$

Furthermore, we note that $(\mathbb{R}^d)^{\otimes m}$ inherits an inner product from \mathbb{R}^d through a choice of basis, as is shown in Equation (3).

Free Algebras. Our goal is to find a vector space H that contains the vector space \mathbb{R}^d but is large enough to support a richer algebraic structure; in particular, a multiplication. Formally, this means we look for an injective map $\mathbb{R}^d \hookrightarrow H$ into an algebra H . A classic mathematical construction that turns a vector space into an algebra is the *free associative algebra*, defined as

$$\bigoplus_{m=0}^{\infty} (\mathbb{R}^d)^{\otimes m} = \{(\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_M, 0, \dots) : \mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}, M \in \mathbb{N}\},$$

where addition and multiplication of two elements $\mathbf{u} = (\mathbf{u}_0, \mathbf{u}_1, \dots)$ and $\mathbf{v} = (\mathbf{v}_0, \mathbf{v}_1, \dots)$ is given by defining the $(\mathbb{R}^d)^{\otimes m}$ coordinate to be

$$(\mathbf{u} + \mathbf{v})_m = \mathbf{u}_m + \mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}, \quad (\mathbf{u} \cdot \mathbf{v})_m = \sum_{i=0}^m \mathbf{v}_i \otimes \mathbf{u}_{m-i} \in (\mathbb{R}^d)^{\otimes m}.$$

We emphasize that the multiplication is not commutative.

A Universal Property. Indeed, this construction is the most general way to turn \mathbb{R}^d into an algebra as the following classical result shows: if we denote with $\iota : \mathbb{R}^d \hookrightarrow \bigoplus_{m=0}^{\infty} (\mathbb{R}^d)^{\otimes m}$ the embedding $\iota(x) = (0, x, 0, \dots)$ then any linear map from \mathbb{R}^d into any associative algebra A factors through ι . In other words, given an associative algebra A and any linear map $h : \mathbb{R}^d \rightarrow A$, there exists a homomorphism of algebras $\tilde{h} : \bigoplus_{m=0}^{\infty} (\mathbb{R}^d)^{\otimes m} \rightarrow H$ such the following diagram commutes

$$\begin{array}{ccc} \mathbb{R}^d & \xrightarrow{h} & H \\ \downarrow \iota & \nearrow \tilde{h} & \\ \bigoplus_{m=0}^{\infty} (\mathbb{R}^d)^{\otimes m} & & . \end{array}$$

In a sense, this shows that $\bigoplus_{m=0}^{\infty} (\mathbb{R}^d)^{\otimes m}$ is the “most general algebra” that \mathbb{R}^d embeds into.

An Inner Product. While the free associative algebra satisfies the above universal property, it is convenient to work in the slightly larger space

$$\prod_{m \geq 0} (\mathbb{R}^d)^{\otimes m} = \{(\mathbf{v}_0, \mathbf{v}_1, \mathbf{v}_2, \dots) : \mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}\}$$

and define an inner product

$$\langle \mathbf{u}, \mathbf{v} \rangle := \sum_{m=0}^{\infty} \langle \mathbf{u}_m, \mathbf{v}_m \rangle_m \quad (16)$$

for elements $\mathbf{u}, \mathbf{v} \in \prod_{m \geq 0} (\mathbb{R}^d)^{\otimes m}$ for which the sum (16) converges (it does not converge for general elements of $\prod_{m \geq 0} (\mathbb{R}^d)^{\otimes m}$). The largest space for which this inner product exists, is

$$H := \left\{ \mathbf{v} \in \prod_{m \geq 0} (\mathbb{R}^d)^{\otimes m} : \|\mathbf{v}\| < \infty \right\} \text{ and } \|\mathbf{v}\| := \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}.$$

The space H contains all the properties we require; it contains \mathbb{R}^d as a subspace, $\iota(\mathbb{R}^d) \subset H$; the non-commutative multiplication structure, $\mathbf{u} \cdot \mathbf{v} \in H$ for $\mathbf{u}, \mathbf{v} \in H$; and a Hilbert space structure, $\langle \mathbf{u}, \mathbf{v} \rangle \in \mathbb{R}$ for $\mathbf{u}, \mathbf{v} \in H$.

Linear Functionals. Any $\ell \in H$ can be treated as a linear functional by taking $\langle \ell, \cdot \rangle : H \rightarrow \mathbb{R}$, we will primarily consider *finite linear functionals*

$$\ell = (\ell_0, \ell_1, \dots, \ell_M, 0, \dots) \in H,$$

which have only finitely many nonzero coordinates. Such finite linear functionals applied to our feature map $\tilde{\varphi}(\mathbf{x}) \in H$ for sequences $\mathbf{x} \in \text{Seq}(\mathbb{R}^d)$, are rich enough so that

$$\mathbf{x} \mapsto \langle \ell, \tilde{\varphi}(\mathbf{x}) \rangle$$

can approximate any functions $f(\mathbf{x})$ of sequences $\mathbf{x} \in \text{Seq}(\mathbb{R}^d)$, and their collection characterizes the distribution of random sequences (i.e. random walks),

$$\mu \mapsto \{ \langle \ell, \mathbb{E}_{\mathbf{x} \sim \mu} [\tilde{\varphi}(\mathbf{x})] \rangle : \ell = (\ell_0, \dots, \ell_M, 0, \dots) \in H, M \geq 1 \}$$

is injective; see Theorem 4.

C The Sequence Feature Map

In Section 2 we introduce the sequence feature map

$$\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H, \quad \tilde{\varphi}(\mathbf{x}) = \varphi(\delta_0 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}),$$

for sequences in \mathbb{R}^d where

$$\varphi : \mathbb{R}^d \rightarrow H$$

is an injective map from \mathbb{R}^d into the algebra H . Here, $\delta_0 \mathbf{x} := x_0$ and $\delta_i \mathbf{x} := x_i - x_{i-1}$ for $i \geq 1$ denote the increments of a sequence $\mathbf{x} = (x_0, \dots, x_k)$. Our main example is the tensor exponential

$$\varphi(x) = \exp_{\otimes}(x) = \left(\frac{x^{\otimes m}}{m!} \right)_{m \geq 0}.$$

The algebra H is the free algebra discussed in detail in Appendix B. In this case, a direct calculation shows that

$$\tilde{\varphi}(\mathbf{x}) = \left(\sum c(i_1, \dots, i_m) \delta_{i_1} \mathbf{x} \otimes \dots \otimes \delta_{i_m} \mathbf{x} \right)_{m \geq 0}$$

where the sum is taken over $i_1 \leq \dots \leq i_m$ with $i_1, \dots, i_m \in \{0, \dots, k\}$ and the coefficients $c(i_1, \dots, i_m) \in \mathbb{R}$ can be computed in explicit form, see [62].

Universality and Characteristicness Universality and characteristicness follow for our “discrete time/sequence” signatures $\tilde{\varphi}(\mathbf{x}) \in H$ from elementary arguments that we discuss below. In our setting it is convenient to include the sequence coordinate in the lifting, that is we consider sequences

$$\bar{\mathbf{x}} = (\bar{x}_0, \dots, \bar{x}_k), \text{ where } \bar{x}_i = (i, x_i) \in \mathbb{R}^{d+1}. \quad (17)$$

Theorem 4. *The time-parametrized sequence feature map*

$$\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H, \quad \tilde{\varphi}(\bar{\mathbf{x}}) = \varphi(\delta_0 \bar{\mathbf{x}}) \cdots \varphi(\delta_k \bar{\mathbf{x}}) \quad (18)$$

has the following properties:

1. **Universality:** *for any continuous⁷ $f : \text{Seq}(\mathbb{R}^d) \rightarrow \mathbb{R}$, any $\epsilon > 0$, and any compact set of sequences $K \subset \text{Seq}(\mathbb{R}^d)$, there exists a $\ell = (\ell_0, \ell_1, \dots, \ell_m, 0, 0, \dots)$ such that*

$$\sup_{\mathbf{x} \in K} |f(\bar{\mathbf{x}}) - \langle \ell, \tilde{\varphi}(\bar{\mathbf{x}}) \rangle| < \epsilon.$$

2. **Characteristicness:** *let $\text{Prob}(K)$ be the set of probability measures that are supported on a compact subset $K \subset \text{Seq}(\mathbb{R}^d)$. Then the map*

$$\text{Prob}(K) \rightarrow H, \quad \mu \mapsto \mathbb{E}_{\mathbf{x} \sim \mu} [\tilde{\varphi}(\bar{\mathbf{x}})]$$

is injective.

Proof. This is a folk theorem in control and probability theory: to see universality, it is sufficient to verify that $\{\langle \ell, \tilde{\varphi}(\bar{\mathbf{x}}) \rangle : \ell = (\ell_0, \dots, \ell_m, 0, \dots), m \geq 0\}$ is a point-separating algebra for $C(K, \mathbb{R})$ since then the result follows from the Stone-Weierstrass theorem. Point separation follows from [22, Corollary 4.9], and a direct calculation shows the product of $\langle \ell, \tilde{\varphi}(\bar{\mathbf{x}}) \rangle$ and $\langle \ell', \tilde{\varphi}(\bar{\mathbf{x}}) \rangle$ is again a linear functional of $\tilde{\varphi}(\bar{\mathbf{x}})$. This shows that the collection of functionals forms an algebra.

The characteristicness follows since $\text{Prob}(K)$ is contained in the dual space of the space of continuous functions of sequences $C(K, \mathbb{R})$ and by universality, linear functionals of $\tilde{\varphi}(\bar{\mathbf{x}})$ are dense in this space, so the result follows. \square

This result shows that:

- non-linear functions of sequences can be approximated using linear functionals in H ; and
- the distribution of random sequences is characterized as the mean of our feature map.

In the terminology of statistical learning this says that $\tilde{\varphi}$ is a universal and characteristic feature map for the set $\text{Seq}(V)$ of sequences. For more background and extensions to non-compact sets of sequences or paths, we refer to [34, 16] and [6, Section 3.2]; for a more geometric picture see [39].

⁷We use the metric $d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^k \|(x_i - y_i) - (x_{i-1} - y_{i-1})\|$ (if two sequences are of different lengths, we pad the end of the shorter one with the end point) to define the topology on $\text{Seq}(\mathbb{R}^d)$.

Path Signatures. We now explain the remark made at the end of Section 2, that for the choice of φ as the tensor exponential (5), the resulting sequence feature map $\tilde{\varphi}$ can be identified as the time discretization of a classical object in analysis, called the path signature. First, we lift a sequence $\mathbf{x} = (x_0, \dots, x_k)$ from discrete time to continuous time by identifying it as the piecewise linear path $\mathbf{X} = (\mathbf{X}(t))_{t \geq 0}$,

$$\mathbf{X}(t) := x_i + (t - i)(x_{i+1} - x_i) \text{ for } t \in [i, i + 1).$$

Then a direct calculation shows that

$$\tilde{\varphi}(\mathbf{x}) = \left(\int_0^k d\mathbf{X}^{\otimes m} \right)_{m \geq 0} \quad \text{where} \quad \int_0^i d\mathbf{X}^{\otimes m} := \int_{0 \leq t_1 \leq \dots \leq t_m \leq i} \dot{\mathbf{X}}(t_1) \otimes \dots \otimes \dot{\mathbf{X}}(t_m) dt_1 \dots dt_m$$

and $i \in [0, k]$; note that $\dot{\mathbf{X}}(t)$ is well-defined for almost every $t \in [0, k]$ since \mathbf{X} is piecewise linear, which is sufficient to make sense of this integral as a Riemann-Lebesgue integral (and stochastic integrals allow to treat rougher paths). Such sequences of iterated integrals are classical in analysis, probability theory, and control theory, and are known under various names (Path-ordered Exponential, Volterra series, Chen–Fliess Series, Chronological Exponential, etc.); we refer to them as *path signatures* as they are known in probability theory.

Diffusions and their Generator. The path signature can even be defined for highly irregular paths such as Brownian motion by using stochastic (Ito-Stratonovich) integrals. This is the connection to our main motivation: if $\mathbf{X} = (\mathbf{X}_t)_{t \geq 0}$ is a Brownian motion in \mathbb{R}^d , then its generator is the classical Laplacian and the diffusion of Brownian particles is captured with the classical diffusion PDE, the heat equation. Gaveau [23] showed that if we lift Brownian trajectories $t \mapsto \mathbf{X}(t)$ evolving in the state space \mathbb{R}^d into paths evolving in a richer state space H via the above signature construction,

$$t \mapsto \left(\int_0^t d\mathbf{X}^{\otimes m} \right)_{m \geq 0}, \quad (19)$$

then the stochastic process (19) is again a Markov process⁸ and its generator satisfies a property called hypo-ellipticity (see the paragraph below). Our approach can be seen as the analogous construction on a graph and in discrete time: in the same way that (19) lifts Brownian motion from \mathbb{R}^d to a process evolving in the state space H , the map $\tilde{\varphi}$ lifts a simple random walk $(B_k)_{k \geq 0}$ on the graph to a random walk $\mathbf{L} = (L_k)_{k \geq 0}$ in H .

Hypo-elliptic Operators, Sub-elliptic Operators, and their Geometry. A full discussion is beyond the scope of this article and several books [31] have been written about this topic. For the interested reader, we give a very short informal picture on the space \mathbb{R}^d which motivates our nomenclature but otherwise this paragraph can be safely skipped. A differential operator $\mathcal{L} = \sum \sigma_{i,i} \frac{\partial^2}{\partial x_i \partial x_j} + \sum_i \mu_i \frac{\partial}{\partial x_j}$ is called elliptic if the matrix $(\sigma_{i,j}(x))_{i,j}$ is invertible for all $x \in \mathbb{R}^d$. Every (time-homogenous) Markov process $\mathbf{X} = (\mathbf{X}_t)_{t \geq 0}$ gives rise to a differential operator

$$\mathcal{L}f(x) = \lim_{t \rightarrow 0} \frac{\mathbb{E}[f(\mathbf{X}_t) | \mathbf{X}_0 = x] - f(x)}{t}$$

but generically \mathcal{L} is not elliptic. An important example class of Markov processes are stochastic differential equations,

$$d\mathbf{X}_t = \sum_{i=1}^e V_i(\mathbf{X}_t) d\mathbf{B}_t^i, \quad \mathbf{X}_0 \in \mathbb{R}^d, \quad (20)$$

where $\mathbf{B}_t = (\mathbf{B}_t^1, \dots, \mathbf{B}_t^e)$ denotes a Brownian motion in \mathbb{R}^e and V_1, \dots, V_e are vector fields on \mathbb{R}^d . By identifying the vectors fields V_1, \dots, V_e as differential operators, the generator of (20) can be written as

$$\mathcal{L} = \sum_{i=1}^e V_i^2, \quad (21)$$

⁸The process (19) is often called the canonical lift of Brownian motion to the free Lie group with d generators. Strictly speaking, Gaveau uses the first M iterated integrals, otherwise one deals with an "infinite-dimensional Lie group" which poses some technical challenges. This Lie group embeds into H and in this article we only use the structure of H and not the Lie group structure; see [23].

which is in general not elliptic. Sub- and hypo-ellipticity are properties of \mathcal{L} that are weaker than ellipticity: ellipticity implies sub-ellipticity and, by a classic theorem of Hörmander, sub-ellipticity implies hypo-ellipticity. Hypo-ellipticity is in turn general enough to cover many important examples for which (sub-)ellipticity fails. For example, another celebrated result of Hörmander is that the SDE generator (21) is hypo-elliptic, whenever the Lie algebra generated by the vector fields and their brackets, spans at every point the full space \mathbb{R}^d . This not only allows us to study the properties of a large class of diffusions but also provides a natural link with geometry: the set of vector fields V_1, \dots, V_e determines a subset of \mathbb{R}^d that the SDE evolves on. For general hypo-elliptic operators, this set is not all of \mathbb{R}^d and it is not even a Riemannian manifold; rather, it is a sub-Riemannian manifold. Intuitively, sub-Riemannian geometry is "rougher" than Riemannian geometry (e.g. no canonical connections exist) but still regular enough so that one can use geometric tools to study the underlying stochastic process and vice versa; see for example [59, 23]. This includes SDEs such as those evolving in free objects: for a natural choice of vector fields V_1, \dots, V_e the solution of the SDE (20) is (19), see [23].

D Further Details on Hypo-elliptic Graph Diffusion

In Section 3, we introduced tensor analogues of classical matrix operators, which were then used to define hypo-elliptic graph diffusion. These tensor-valued matrices allow us to efficiently represent random walk histories on a graph, which can be manipulated using matrix operations. In this section, we discuss further details on the tensor adjacency matrix, provide a proof of Theorem 1, and introduce a variation of hypo-elliptic graph diffusion.

As in Section 3, we fix a labelled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, f)$, where the nodes are denoted by integers, $\mathcal{V} = \{1, \dots, n\}$, and $f : \mathcal{V} \rightarrow \mathbb{R}^d$ denotes the continuous node attributes.

Powers of the Tensor Adjacency Matrix. Recall that the powers of the classical adjacency matrix A counts the number of walks between two nodes on the graph. In particular, given $k \in \mathbb{N}$, and two nodes $i, j \in \mathcal{V}$, the result follows as a consequence of the sparsity pattern of A ,

$$(A^k)_{i,j} = \sum_{i_1, \dots, i_{k-1}=1}^n A_{i,i_1} \cdot A_{i_1,i_2} \cdots A_{i_{k-1},j} = \sum_{i=i_0 \sim \dots \sim i_k=j} 1.$$

Note that the product $A_{i,i_1} \cdots A_{i_{k-1},j} = 0$ unless each pair of consecutive indices are adjacent in the graph, namely $i_{q-1} \sim i_q$ for all $q = 1, \dots, k$. Applying the same procedure to the tensor adjacency matrix from Equation (10), we obtain a summary of all walks between two nodes, rather than just the number of walks. In particular,

$$\begin{aligned} (\tilde{A}^k)_{i,j} &= \sum_{i=i_0 \sim \dots \sim i_k=j} \tilde{A}_{i,i_1} \cdot \tilde{A}_{i_1,i_2} \cdots \tilde{A}_{i_{k-1},j} \\ &= \sum_{i=i_0 \sim \dots \sim i_k=j} \varphi(f(i_1) - f(i)) \cdot \varphi(f(i_2) - f(i_1)) \cdots \varphi(f(j) - f(i_{k-1})) \\ &= \sum_{i=i_0 \sim \dots \sim i_k=j} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}), \end{aligned}$$

where $\mathbf{x} = (f(i_0), \dots, f(i_k))$ denotes the lifted sequence in the vector space \mathbb{R}^d . Note that this corresponds to the sequence feature map *without* the initial point $\delta_0 \mathbf{x}$.

Powers of the Tensor Transition Matrix. We now consider powers of the classical transition matrix $P = I - \mathcal{L}$, where $\mathcal{L} = I - D^{-1}A$ is the normalized graph Laplacian. The entries of P^k provide length k random walk probabilities; in particular, we have

$$(P^k)_{i,j} = \sum_{i=i_0 \sim \dots \sim i_k=j} \frac{1}{\deg(i_0) \deg(i_1) \cdots \deg(i_{k-1})} = \mathbb{P}[B_k = j | B_0 = i].$$

The powers of the tensor transition matrix $\tilde{P} = I - \tilde{\mathcal{L}}$, where $\tilde{\mathcal{L}} = I - D^{-1}\tilde{A}$ is the hypo-elliptic Laplacian, will be the conditional expectation of the sequence feature map of the random walk

process. In particular,

$$(\tilde{P})_{i,j}^k = \sum_{i=i_0 \sim \dots \sim i_k=j} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}) \mathbb{P}[B_1 = i_1, B_2 = i_2, \dots, B_k = i_k | B_0 = i]. \quad (22)$$

Proof of Hypo-elliptic Diffusion Theorem. Recall from Section 3 that the *hypo-elliptic graph diffusion equation* is defined by

$$\mathbf{v}_k - \mathbf{v}_{k-1} = -\tilde{\mathcal{L}}\mathbf{v}_{k-1}, \quad \mathbf{v}_0 = \mathbb{1}_H,$$

where $\mathbb{1}_H^T = (1_H, \dots, 1_H) \in H^n$ is the all-ones vector in H . Using the above computations for powers of the tensor transition matrix, we can prove Theorem 1, which is restated here.

Theorem 1. *Let $k \in \mathbb{N}$, $\mathbf{L}_k = (L_0, \dots, L_k)$ be the lifted random walk from (6), and $\tilde{P} = I - \tilde{\mathcal{L}}$ be the tensor adjacency matrix. The solution to the hypo-elliptic graph diffusion equation (12) is*

$$\mathbf{v}_k = (\mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_0 = i])_{i=1}^n = \tilde{P}^k \mathbb{1}_H.$$

Furthermore, if $F \in H^{n \times n}$ is the diagonal matrix with $F_{i,i} = \varphi(f(i))$, then

$$F\mathbf{v}_k = (\mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i])_{i=1}^n.$$

Proof of Theorem 1. We begin by proving the first equation. From the definition of hypo-elliptic diffusion, it is straightforward to see that

$$\mathbf{v}_k = (I - \tilde{\mathcal{L}})\mathbf{v}_{k-1} = \tilde{P}\mathbf{v}_{k-1} = \tilde{P}^k \mathbf{v}_0 = \tilde{P}^k \mathbb{1}_H.$$

We prove coordinate-wise that $\mathbf{v}_k = (\mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_0 = i])_{i=1}^n$. Indeed, using the above and Equation (22), the i -th coordinate of \mathbf{v}_k is

$$\begin{aligned} \mathbf{v}_k^{(i)} &= (\tilde{P}^k \mathbb{1}_H)^{(i)} = \sum_{j=1}^n (\tilde{P}^k)_{i,j} = \sum_{i=i_0 \sim \dots \sim i_k} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}) \mathbb{P}[B_1 = i_1, \dots, B_k = i_k | B_0 = i] \\ &= \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_0 = i], \end{aligned}$$

where $\mathbf{x} = (f(i_0), \dots, f(i_k))$ is the lifted sequence corresponding to a walk $i_0 \sim \dots \sim i_k$ on the graph. Next, we will also prove the second equation coordinate-wise. Using the above result, we have

$$(F\mathbf{v}_k)^{(i)} = \varphi(f(i)) \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_0 = i] = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i],$$

where we use the fact that $\delta_0 \mathbf{L}_k = L_0 = f(i)$ when we condition $B_0 = i$. \square

Forward Hypo-elliptic Diffusion. In the classical setting, we can consider both the forward and backward Kolmogorov equations. Throughout the main text and in the appendix so far, we have been considering the *backward* variants. In this section, we formulate the *forward* analogue of hypo-elliptic diffusion. In the classical graph setting, this corresponds to the following forward equation for $U_k \in \mathbb{R}^{n \times d}$ given by

$$U_k^T - U_{k-1}^T = -U_{k-1}^T \mathcal{L}, \quad U_0^{(i)} = f(i)$$

where the initial condition $U_0 \in \mathbb{R}^{n \times d}$ is specified by the node attributes. Note that because $P = D^{-1}A$ is right stochastic⁹, this variation of the graph diffusion equation conserves mass in each coordinate of the node attributes at every time step.

Similarly we can formulate the forward hypo-elliptic graph diffusion equation for $\mathbf{v}_k \in H^n$ as

$$\mathbf{v}_k^T - \mathbf{v}_{k-1}^T = -\mathbf{v}_{k-1}^T \tilde{\mathcal{L}}, \quad \mathbf{v}_0 = n^{-1} \mathbb{1}_H. \quad (23)$$

The solution of the Equation 23 is given below.

⁹Each column sums to one and hence multiplying on the left by a row vector conserves its mass.

Theorem 5. Let $k \in \mathbb{N}$, $\mathbf{L}_k = (L_0, \dots, L_k)$ be the lifted random walk from (6), and $\tilde{P} = I - \tilde{\mathcal{L}}$ be the tensor adjacency matrix. The solution to the forward hypo-elliptic graph diffusion equation (23) is

$$\mathbf{v}_k^T = (\mathbb{P}[B_k = i] \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_k = i])_{i=1}^n = \frac{1}{n} \mathbb{1}_H^T \tilde{P}^k.$$

Furthermore, if $F \in H^{n \times n}$ is the diagonal matrix with $F_{i,i} = \varphi(f(i))$, then

$$\frac{1}{n} \mathbb{1}_H^T F \tilde{P}^k = (\mathbb{P}[B_k = i] \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_k = i])_{i=1}^n.$$

Proof. The proof proceeds in the same way as the backward equation. By definition of the forward hypo-elliptic diffusion, we have

$$\mathbf{v}_k^T = \mathbf{v}_0^T \tilde{P}^k = \frac{1}{n} \mathbb{1}_H^T \tilde{P}^k.$$

Now, recall that the initial point of the random walk process is chosen uniformly over all nodes; in other words, $\mathbb{P}[B_0 = i] = \frac{1}{n}$. Then, we show $\mathbf{v}_k^T = (\mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_k = i])_{i=1}^n$ coordinate-wise as

$$\begin{aligned} \mathbf{v}_k^{(i)} &= \frac{1}{n} \sum_{j=1}^n (\tilde{P}^k)_{j,i} \\ &= \sum_{i_0 \sim \dots \sim i_k = i} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}) \mathbb{P}[B_1 = i_1, \dots, B_k = i | B_0 = i_0] \mathbb{P}[B_0 = i_0] \\ &= \sum_{i_0 \sim \dots \sim i_k = i} \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}) \mathbb{P}[B_0 = j, \dots, B_{k-1} = i_{k-1} | B_k = i] \mathbb{P}[B_k = i] \\ &= \mathbb{P}[B_k = i] \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) | B_k = i], \end{aligned}$$

where $\mathbf{x} = (f(i_0), \dots, f(i_k))$ is the lifted sequence corresponding to a walk $i_0 \sim \dots \sim i_k$ on the graph. We will now prove the second equation. Note that we have

$$(F \tilde{P}^k)_{i,j} = \sum_{i=i_0 \sim \dots \sim i_k = j} \tilde{\varphi}(\mathbf{x}) \mathbb{P}[B_1 = i_1, \dots, B_k = i_k | B_0 = i].$$

Then, following the same reasoning as the first equation, we obtain the desired result. \square

Weighted Graphs. In the prior discussion, we considered simple random walks in which the walk chooses one of the nodes adjacent to the current node uniformly at random. We can instead consider more general random walks on graphs, which can be described using a weighted adjacency matrix,

$$A_{i,j} = \begin{cases} c_{i,j} & : i \sim j \\ 0 & : \text{otherwise.} \end{cases}$$

In this case, we define the diagonal weighted degree matrix to be $D_{i,i} = \sum_{i \sim j} A_{i,j}$. We now define a weighted random walk $(B_k)_{k \geq 0}$ on the vertices \mathcal{V} where the transition matrix is given by

$$P_{i,j} := D^{-1} A = \mathbb{P}(B_k = j | B_{k-1} = i) = \begin{cases} \frac{c_{i,j}}{\sum_{i \sim j'} c_{i,j'}} & : i \sim j \\ 0 & : \text{otherwise.} \end{cases}$$

With these weighted graphs, the powers of the adjacency and transition matrix can be interpreted in a similar manner as the standard case. Powers of the adjacency matrix provide total weights over walks, while the powers of the transition matrix provide the probability of a weighted walk going between two nodes after a specified number of steps. In particular,

$$\begin{aligned} (A^k)_{i,j} &= \sum_{i \sim i_1 \sim \dots \sim i_{k-1} \sim j} c_{i,i_1} \cdots c_{i_{k-1},j} \\ (P^k)_{i,j} &= \mathbb{P}[B_k = j | B_0 = i]. \end{aligned}$$

The tensor adjacency and tensor transition matrices are defined in the same manner as

$$\begin{aligned} \tilde{A}_{i,j} &:= \begin{cases} c_{i,j} \varphi(f(j) - f(i)) & : i \sim j \\ 0 & : \text{otherwise} \end{cases} \\ \tilde{P}_{i,j} &:= D^{-1} \tilde{A} = \begin{cases} P_{i,j} \varphi(f(j) - f(i)) & : i \sim j \\ 0 & : \text{otherwise} \end{cases} \end{aligned}$$

Note that powers of this weighted tensor transition matrix are exactly the same as the unweighted case from Equation (22), and thus the weighted version of both Theorem 2 and Theorem 5 immediately follow.

E Characterizing Random Walks

In this appendix, we will provide further details on how the features obtained via hypo-elliptic diffusion are able to characterize the underlying random walk processes. These results rely on the characteristic property of the time-parametrized sequence feature map from Theorem 4.

Computation with Time Parametrization. Recall that the time-parametrized sequence feature map from Equation (18) applies the algebra lifting to the time-parametrized sequence $\bar{\mathbf{x}} := (\bar{x}_0, \dots, \bar{x}_k)$, where $\bar{x}_i = (i, x_i)$. Note that we have $\delta_i \bar{\mathbf{x}} = (1, x_i - x_{i-1})$. Thus, the hypo-elliptic diffusion equations and the low-rank algorithm given in Theorem 3 easily extends to this setting.

Characterizing Random Walks. Recall that hypo-elliptic diffusion yields a feature map for labelled graphs by mean-pooling the individual node features as

$$\Psi(\mathcal{G}) = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k)],$$

where $\mathbf{L}_k = (L_0, \dots, L_k)$ is the lifted random walk process in \mathbb{R}^d . We will now prove Theorem 2, which we restate here with more details.

Theorem 2. *Let*

$$\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H, \quad \tilde{\varphi}(\bar{\mathbf{x}}) = \varphi(\delta_0 \bar{\mathbf{x}}) \cdots \varphi(\delta_k \bar{\mathbf{x}}),$$

where $\bar{\mathbf{x}}$ appends the time parametrization to the sequence \mathbf{x} as in Equation (17). Furthermore, suppose we have the resulting graph feature map

$$\Psi(\mathcal{G}) = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k)].$$

Let \mathcal{G} and \mathcal{G}' be two labelled graphs, and $\mathbf{L}_k = (L_0, \dots, L_k)$ and $\mathbf{L}'_k = (L'_0, \dots, L'_k)$ be the k -step lifted random walk as defined in Equation (6), given the random walk processes B and B' on \mathcal{G} and \mathcal{G}' respectively. Then, $\Psi(\mathcal{G}) = \Psi(\mathcal{G}')$ if and only if the distributions of \mathbf{L}_k and \mathbf{L}'_k are equal.

Proof. First, if the distributions of the two random walks \mathbf{L}_k and \mathbf{L}'_k are equal, then it is clear that $\mathbb{E}[\tilde{\varphi}(\mathbf{L}_k)] = \mathbb{E}[\tilde{\varphi}(\mathbf{L}'_k)]$.

Next, suppose $\mathbb{E}[\tilde{\varphi}(\mathbf{L}_k)] = \mathbb{E}[\tilde{\varphi}(\mathbf{L}'_k)]$. Then, note that the random walk distributions are finitely supported (hence compactly supported) distributions, and thus by Theorem 4, they must be equal. \square

This result shows that the hypo-elliptic diffusion completely characterizes random walk histories; thus providing a highly descriptive summary of labelled graphs. There is an analogous result for the individual node features in terms of conditional expectations.

Theorem 6. *Let*

$$\tilde{\varphi} : \text{Seq}(\mathbb{R}^d) \rightarrow H, \quad \tilde{\varphi}(\bar{\mathbf{x}}) = \varphi(\delta_0 \bar{\mathbf{x}}) \cdots \varphi(\delta_k \bar{\mathbf{x}}),$$

where $\bar{\mathbf{x}}$ appends the time parametrization to the sequence \mathbf{x} as in Equation (17). Furthermore, suppose we have the resulting node feature map

$$\Phi(i) = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = i].$$

Let \mathcal{G} and \mathcal{G}' be two labelled graphs, and $\mathbf{L}_k = (L_0, \dots, L_k)$ and $\mathbf{L}'_k = (L'_0, \dots, L'_k)$ be the k -step lifted random walk as defined in Equation (6), given the random walk processes B and B' on \mathcal{G} and \mathcal{G}' respectively. Then for two nodes $i \in \mathcal{V}$ and $i' \in \mathcal{V}'$ on the two respective graphs, $\Phi(i) = \Phi(i')$ if and only if the conditional distributions of $\mathbb{P}[\mathbf{L}_k | B_0 = i]$ and $\mathbb{P}[\mathbf{L}'_k | B'_0 = i']$ are equal.

Proof. The proof is analogous to the proof of Theorem 2 given above. \square

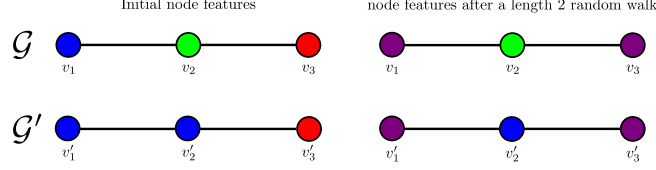


Figure 1: On the left: the graphs \mathcal{G} and \mathcal{G}' with their initial labels. On the right: the labels of \mathcal{G} and \mathcal{G}' after random walk of length 2.

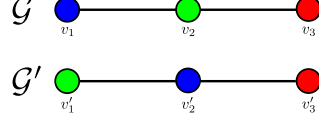


Figure 2: Here the global graph features computed using classical diffusion cannot distinguish between the graphs \mathcal{G} and \mathcal{G}' , however using our framework creates global graph features that can distinguish between the two.

Examples. As mentioned after Theorem 2 in the main text, the benefits of capturing path-dependence can be significant. In addition to the experimental results one can gain intuition for this by looking at simple toy examples. For example, consider as labels RGB colors, that is $f : \mathcal{V} \rightarrow [0, 1]^3 \subset \mathbb{R}^3$ and random walks of length $k = 2$ on the two graphs $\mathcal{G}, \mathcal{G}'$ with vertex sets $\mathcal{V} = \{v_1, v_2, v_3\}$, resp. $\mathcal{V}' = \{v'_1, v'_2, v'_3\}$.

We will now see two examples where features computed using classical diffusion cannot distinguish either the nodes or the graphs, but features computed using hypo-elliptic diffusion can.

Consider colors $a, b, c \in \mathbb{R}^3$ as labels on the nodes of a graph with three nodes and two edges as are represented in Figures 1 and 2, with colors a and c on the end nodes and b on the middle node. After a random walk B_2 of length two, simple calculations show that the labels on the nodes will be $\frac{1}{2}(a + c)$ on the two end nodes, and b on the middle node.

The first example in Figure 1 considers node features. If one only has access to the marginal distribution $L_2 = f(B_2)$ to construct node features, then the left-most and right-most nodes in the graphs \mathcal{G} and \mathcal{G}' from Figure 1 are indistinguishable, that is

$$\mathbb{E}[L_2 | B_0 = v_1] = \mathbb{E}[L'_2 | B'_0 = v'_1] = \frac{1}{2}(r + b).$$

In stark contrast, the laws of (L_0, L_1, L_2) and (L'_0, L'_1, L'_2) clearly differ since L_1 , resp. L'_1 , are with probability one equal to green, resp. blue. Thus Theorem 2 guarantees that

$$\Phi(v_1) = \mathbb{E}[\tilde{\varphi}(\mathbf{L}_k) | B_0 = v_1] \neq \mathbb{E}[\tilde{\varphi}(\mathbf{L}'_k) | B'_0 = v'_1] = \Phi(v'_1).$$

In fact, for this example, a direct calculation shows that

$$\langle e_b \otimes e_g, \tilde{\varphi}(\mathbf{L}_2) \rangle = 1 \neq 0 = \langle e_b \otimes e_g, \tilde{\varphi}(\mathbf{L}'_2) \rangle.$$

The second example in Figure 2 considers graph features. The global graph features for 2-step classical diffusion is computed as the average of the node features, i.e. the global graph feature is $\frac{1}{3}(a + b + c)$. Hence, the graphs of Figure 2 are not distinguishable through features computed from standard diffusion. However, using hypo-elliptic diffusion allows us to take into account the order in which the colors appear in the random walk. Indeed the colors red and green never appear consecutively in the graph \mathcal{G}' , hence the $e_r \otimes e_g$ component of $\Psi(\mathcal{G}')$ is zero; in other words,

$$\langle e_r \otimes e_g, \Psi(\mathcal{G}) \rangle = 1 \neq 0 = \langle e_r \otimes e_g, \Psi(\mathcal{G}') \rangle,$$

showing that $\Psi(\mathcal{G}) \neq \Psi(\mathcal{G}')$.

Simple examples as the above demonstrate how using only the marginal distribution after k steps lacks expressive power for both node and graph features. Therefore, while classical diffusion smooths out neighbourhood features, hypo-elliptic diffusion retains finer descriptions of local neighbourhoods.

F Details on the Low Rank Algorithm.

In this section, we will provide further details and proofs on the low-rank approximation method discussed in Section 4. We can use low-rank tensors to define corresponding low-rank functionals of the features obtained via hypo-elliptic diffusion. The following is the definition of *CP-rank* of a tensor from [10].

Definition 1. The *rank* of a level m tensor $\mathbf{v}_m \in (\mathbb{R}^d)^{\otimes m}$ is the smallest number $r \geq 0$ such that we can express \mathbf{v}_m as

$$\mathbf{v}_m = \sum_{i=1}^r v_{i,1} \otimes \dots \otimes v_{i,m}, \quad v_{i,j} \in \mathbb{R}^d.$$

We say that $\mathbf{v} = (\mathbf{v}_0, \mathbf{v}_1, \dots) \in H$ is *rank 1* if all \mathbf{v}_m are rank 1 tensors.

We will now prove Theorem 3, which is stated using the node feature map without ZeroStart (see Appendix G). In particular, the node feature map is given by

$$\hat{\Phi}_k(i) = \mathbb{E}[\varphi(\delta_1 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) \mid B_0 = i] = (\tilde{P}^k \mathbb{1}_H)^{(i)} \in H, \quad (24)$$

where we explicitly specify the walk length k in the subscript. Furthermore, if we also need to specify the tensor degree m , we will use two subscripts, where

$$\hat{\Phi}_{k,m}(i) \in (\mathbb{R}^d)^{\otimes m}$$

is the level m component of the hypo-elliptic diffusion with a walk length of k . Throughout the proof, we omit the H subscript for the all-ones vector, such that $\mathbb{1}^T := (1_H, \dots, 1_H) \in H^n$, and we denote $\mathbb{1}_i^T = (0, \dots, 1_H, \dots, 0) \in H^n$ to be the unit vector in the i^{th} coordinate.

Proof of Theorem 3. First, we will show that $f_{1,m}(i) = \langle \ell_m, \hat{\Phi}_1(i) \rangle$ for all $m = 1, \dots, M$. By the definition of hypo-elliptic diffusion, we know that

$$\hat{\Phi}_1(i) = \mathbb{1}_i^T \tilde{P} \mathbb{1} = \sum_{j=1} \tilde{P}_{i,j} = \sum_{i \sim j} \frac{\exp_{\otimes}(f(j) - f(i))}{d_i}.$$

By explicitly expressing the level m component, and by factoring out the inner product, we get

$$\begin{aligned} \left\langle \ell_m, \frac{\exp_{\otimes}(f(j) - f(i))}{d_i} \right\rangle &= \frac{1}{d_i m!} \prod_{r=M-m+1}^M \langle u_r, f(j) - f(i) \rangle \\ &= \frac{1}{m!} (P_{i,j} \cdot C_{i,j}^{u_{M-m+1}} \cdots C_{i,j}^{u_M}). \end{aligned}$$

Then by linearity of the inner product, we get $f_{1,m}(i) = \langle \ell_m, \hat{\Phi}_1(i) \rangle$.

Next, we continue by induction and suppose that $f_{k-1,m}(i) = \langle \ell_m, \hat{\Phi}_{k-1}(i) \rangle$ holds for all $m = 1, \dots, M$. Starting once again from the definition of hypo-elliptic diffusion, we know that

$$\hat{\Phi}_k(i) = \mathbb{1}_i^T \tilde{P} \hat{\Phi}_{k-1} = \sum_{i \sim j} \tilde{P}_{i,j} \cdot \hat{\Phi}_{k-1}(j).$$

Fix a degree m . We explicitly write out the level m component of this equation by expanding $\tilde{P}_{i,j}$ and the tensor product as

$$\begin{aligned} \hat{\Phi}_{k,m}(i) &= \sum_{i \sim j} \sum_{r=0}^m \frac{(f(j) - f(i))^{\otimes r}}{d_i r!} \cdot \hat{\Phi}_{k-1,m-r}(j) \\ &= \sum_{i \sim j} \frac{\hat{\Phi}_{k-1,m}(j)}{d_i} + \sum_{r=1}^m \frac{1}{r!} \sum_{i \sim j} \frac{(f(j) - f(i))^{\otimes r}}{d_i} \cdot \hat{\Phi}_{k-1,m-r}(j) \end{aligned} \quad (25)$$

Note that the first sum is equivalent to

$$\sum_{i \sim j} \frac{\hat{\Phi}_{k-1,m}(j)}{d_i} = \mathbb{1}_i^T P \cdot \hat{\Phi}_{k-1,m}.$$

Applying the linear functional ℓ_m and the induction hypothesis to this, we have

$$\left\langle \ell_m, \sum_{i \sim j} \frac{\hat{\Phi}_{k-1,m}(j)}{d_i} \right\rangle = \mathbb{1}_i^T P \cdot f_{k-1,m}.$$

For the second sum in Equation (25), we can factor the inner product and apply the induction hypothesis to get

$$\left\langle \ell_m, \sum_{i \sim j} \frac{(f(j) - f(i))^{\otimes r}}{d_i} \cdot \hat{\Phi}_{k-1,m-r}(j) \right\rangle = \sum_{j=1}^n P_{i,j} \cdot C_{i,j}^{u_{M-m+1}} \cdot \dots \cdot C_{i,j}^{u_{M-m+r}} \cdot f_{k-1,m-r}.$$

Putting this all together, we get

$$\hat{\Phi}_{k,m}(i) = \mathbb{1}_i^T \left(P \cdot f_{k-1,m} + \sum_{r=1}^m \frac{1}{r!} (P \odot C^{u_{M-m+1}} \odot \dots \odot C^{u_M}) \cdot f_{k-1,m-r} \right) = f_{k,m}(i).$$

□

Computing the ZeroStart Variation. We can adapt the recursive algorithm provided by Theorem 3 above in order to compute low rank approximations to the ZeroStart variation of the node features, which is used throughout the main text (see also Appendix G). In particular, we consider

$$\Phi_k(i) = \mathbb{E}[\varphi(\delta_0 \mathbf{L}_k) \cdots \varphi(\delta_k \mathbf{L}_k) \mid B_0 = i] = \varphi(\delta_0 \mathbf{L}_k) \cdot \hat{\Phi}_k(i), \quad (26)$$

where $\hat{\Phi}_k$ is the variation without ZeroStart defined in Equation (24). Note that we can factor the $\varphi(\delta_0 \mathbf{L}_k)$ term out of the expectation due to the conditioning $B_0 = i$, and thus $\varphi(\delta_0 \mathbf{L}_k) = \varphi(f(i))$ is fixed. Thus, we can compute low rank functionals of $\Phi_k(i)$ using one additional step.

Theorem 7. *Using the same hypotheses as Theorem 3; let*

$$\ell_m = u_{M-m+1} \otimes \dots \otimes u_M,$$

where $u_m \in \mathbb{R}^d$ for $m = 1, \dots, M$ and let

$$F_i^u := \langle u, f(i) \rangle.$$

Then,

$$\langle \ell_M, \Phi_k(i) \rangle = \sum_{r=0}^M \frac{1}{m!} F_i^{u_1} \cdots F_i^{u_r} \cdot f_{k,M-r}(i),$$

where $f_{k,m}(i) = \langle \ell_m, \hat{\Phi}_k(i) \rangle$ from Theorem 3.

Proof. Using the definition of $\Phi_k(i)$ from Equation (26), and expanding out the definition of $\varphi(\delta_0 \mathbf{L}_k)$ at level M , we have

$$\Phi_{k,M}(i) = \sum_{r=0}^M \frac{f(i)^{\otimes r}}{r!} \cdot \hat{\Phi}_{k,M-r}(i).$$

Then, taking the linear functional and distributing, we obtain the result

$$\langle \ell_M, \Phi_k(i) \rangle = \sum_{r=0}^M \frac{1}{r!} F_i^{u_1} \cdots F_i^{u_r} \cdot f_{k,M-r}(i).$$

□

Algorithm 1 Computing a rank-1 functional of the hypoelliptic node features.

```

1: Input: Graph  $\mathcal{G} = ([n], \mathcal{E}, f)$  with adjacency matrix  $A$  and transition matrix  $P$ 
2: a rank-1 tensor  $\ell = u_1 \otimes \dots \otimes u_M \in T(\mathbb{R}^d)$ , truncation level  $M \in \mathbb{N}$ , walk length  $k \in \mathbb{N}$ 
3: Compute  $C[m, i, j] \leftarrow A[i, j] \cdot \langle u_m, f(j) - f(i) \rangle$  for  $m \in [M]$ ,  $r \in [R]$ , and  $i, j \in [n]$ 
4: for  $i = 1$  to  $k$  do
5:   parfor  $m = 1$  to  $M$  do
6:     if  $i == 1$  then
7:       Assign  $Q[1, m, :] \leftarrow \frac{1}{m!} (P \odot C[M - m + 1, :, :] \odot \dots \odot C[M, :, :]) \cdot \mathbb{1}$ 
8:     else
9:       Assign  $Q[i, m, :] \leftarrow P \cdot Q[i - 1, m, :] + \sum_{j=1}^m \frac{1}{j!} (P \odot C[M - m + 1, :, :] \odot \dots \odot C[M - m + j, :]) \cdot Q[i - 1, m - j, :]$ 
10:    end if
11:  end parfor
12: end for
13: Output: Vector  $Q[k, m, :]$  of length  $n$  representing node features  $(\langle \ell, \hat{\Phi}_k(i) \rangle)_{i=1, \dots, n}$ 

```

Computational Complexity. We will now consider the computational complexity of our algorithms. We begin by noting that the naive approach of computing $\Phi_k(i) = (F\tilde{P}^k \mathbb{1}_H)^{(i)}$ has the computational complexity of matrix multiplication; though this counts tensor operations, which itself requires $O(d^m)$ scalar multiplications at tensor degree m . This is computationally too expensive for practical applications.

Next, we consider the complexity of the recursive low-rank algorithm from Theorem 3, where the primary computational advantage is the fact that we only perform *scalar* operations rather than *tensor* operations. We consider the recursive step from Equation (15), reproduced here for $m = M$,

$$f_{k,M} := P \cdot f_{k-1,M} + \sum_{r=1}^M \frac{1}{r!} (P \odot C^{u_1} \odot \dots \odot C^{u_r}) \cdot f_{k-1, M-r}.$$

Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, f)$ with $n = |\mathcal{V}|$ nodes and $E = |\mathcal{E}|$ edges, both P and C^u are sparse $n \times n$ (scalar) matrices with $O(E)$ nonzero entries, and $f_{k,m}$ is an n -dimensional column vector. Recall that \odot denotes element-wise multiplication, and thus both the sparse matrix-matrix multiplication and the sparse matrix-vector multiplication have complexity $O(E)$. Furthermore, the entry-wise products $C^{u_1} \odot \dots \odot C^{u_r}$ differ by only one factor between $r = m$ and $r = m + 1$, and thus, computing $f_{k,M}$ assuming all lower $f_{k',m'}$ have been computed has complexity $O(ME)$. Taking into account the two recursive parameters results in a complexity of $O(kM^2E)$. Note that this is the complexity to compute features for *all nodes*.

Once the $f_{k,m}$ are computed, the complexity of adding the start point from Theorem 7 is $O(M)$.

Pseudocode. We provide pseudocode in Algorithm 1 to demonstrate the implementation of Theorem 3. The primary idea is to recursively compute the $f_{k,m}$ vector from Theorem 3 (renamed as an array $Q[:, :, :]$ to avoid confusion with the node features f). In particular, note that the inner loop over m can be parallelized. Thus, while the theoretical complexity in terms of max tensor level is $O(M^2)$ as discussed above, when we run the algorithm in parallel on a GPU, the computation time scales roughly linearly, as shown empirically in Table 7 in Appendix H.4.

G Variations and Hyperparameters of Hypo-Elliptic Diffusion

In this appendix, we summarize possible variations of the sequence feature map, leading to different hypo-elliptic diffusion features. The choice of variation is learned during training, and we also summarize the hyperparameters used for our features. While the theoretical results on characterizing random walks, such as Theorem 2, depend on specific choices of the sequence feature map, there exist analogous results for these variations, which can characterize random walks up to certain equivalences. Furthermore, the computation of these variations can be performed in the same way: through tensorized linear algebra for exact solutions, and through an analogous low-rank method (as in Theorem 3) for approximate solutions.

We fix an algebra lifting $\varphi : \mathbb{R}^d \rightarrow H$ and let $\mathbf{x} = (x_0, \dots, x_k) \in \text{Seq}(\mathbb{R}^d)$. The simplest sequence feature map to define simply multiplies the terms in the sequence together as

$$\tilde{\varphi}(\mathbf{x}) = \varphi(x_0) \cdots \varphi(x_k). \quad (27)$$

Note that this is *not* the sequence feature map used in the main text. We will now discuss several variations of this map, where $\tilde{\varphi}_{\text{inc,zs}}$ is the one primarily used in the main text and $\tilde{\varphi}_{\text{inc,zs,tp}}$ is used to characterize random walks in Theorem 2 and Appendix E.

Increments (Diff). Rather than directly multiplying terms in the sequence together, we can instead multiply the *increments* as

$$\tilde{\varphi}_{\text{inc}}(\mathbf{x}) = \varphi(\delta_1 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}),$$

where $\delta_i \mathbf{x} := x_i - x_{i-1}$ for $i \geq 1$. In both cases, the sequence feature map is the path signature of a continuous piecewise-linear path when we set $\varphi = \exp_{\otimes}$, as discussed in Appendix C, and it is instructive to use this perspective to understand the effect of increments. If we use increments, the path corresponding to the sequence is

$$\mathbf{X}_{\text{inc}}(t) := x_i + (t - i)(x_{i+1} - x_i) \text{ for } t \in [i, i + 1),$$

while if we do not use increments, the path corresponding to the sequence is

$$\mathbf{X}(t) := \sum_{j=0}^{i-1} x_j + (t - i)x_i \text{ for } t \in [i, i + 1).$$

Thus, when we use increments the sequence \mathbf{x} corresponds to the vertices of the path \mathbf{X}_{inc} , while if we do not, it corresponds to the vectors between vertices of the path \mathbf{X} . In practice, this variation corresponds to taking first-differences of the sequence \mathbf{x} before using eq. (27).

Zero starting point (ZeroStart). The sequence feature map with increments, $\tilde{\varphi}_{\text{inc}}$, as defined above is *translation-invariant*, meaning $\tilde{\varphi}_{\text{inc}}(\mathbf{x} + a) = \tilde{\varphi}_{\text{inc}}(\mathbf{x})$, where $\mathbf{x} + a = (x_0 + a, x_1 + a, \dots, x_k + a)$ for some $a \in \mathbb{R}^d$. In order to remove translation invariance, we can start each sequence at the origin $0 \in \mathbb{R}^d$ by pre-appending a 0 to each sequence. A concise way to define the resulting *zero started* sequence feature map is

$$\tilde{\varphi}_{\text{inc,zs}}(\mathbf{x}) = \varphi(\delta_0 \mathbf{x}) \cdots \varphi(\delta_k \mathbf{x}),$$

where we define $\delta_0 \mathbf{x} := x_0$. This is the sequence feature map defined in Equation (1). Note that this variation does not change the sequence feature map if we do not use increments.

Time parametrization. When we relate sequences to piecewise linear paths as described in Appendix C, we can use the fact that the path signature is invariant under reparametrization, or more generally, tree-like equivalence [28]. In terms of discrete sequences, this includes invariance with respect to 0 elements in the sequence (without increments), and repeated elements in the sequence (with increments). In order to remove this invariance, we can include *time parametrization* by setting

$$\tilde{\varphi}_{-, \text{tp}}(\mathbf{x}) := \tilde{\varphi}_-(\bar{\mathbf{x}}),$$

where $\bar{\mathbf{x}} := (\bar{x}_0, \dots, \bar{x}_k) \in \text{Seq}(\mathbb{R}^{d+1})$, with $\bar{x}_i := (i, x_i) \in \mathbb{R}^{d+1}$. This is a simple form of positional encoding, but other encodings are also possible, e.g. sinusoidal waves as in [64].

Algebra lifting (AlgOpt). Throughout this article, we have used the tensor exponential as the algebra lifting. However, we can also scale each level of the lifting independently, and keep these as hyperparameters to optimize. In particular, for a sequence $\mathbf{c} = (c_0, c_1, \dots) \in \mathbb{R}^{\mathbb{N}}$, we define $\varphi^{\mathbf{c}} : \mathbb{R}^d \rightarrow H$ to be

$$\varphi^{\mathbf{c}}(x) := \left(c_m x^{\otimes m} \right)_{m=0}^{\infty},$$

where $1/m!$ is used as initialization for c_m and learned along with the other parameters.

The choice of which variant of the sequence feature map to use depends on which invariance properties are important for the specific problem. In practice, the choice can be learned during the training, which is done in our experiments in Section 5. Furthermore, the features obtained through the low-rank hypo-elliptic diffusion depend on three hyperparameters:

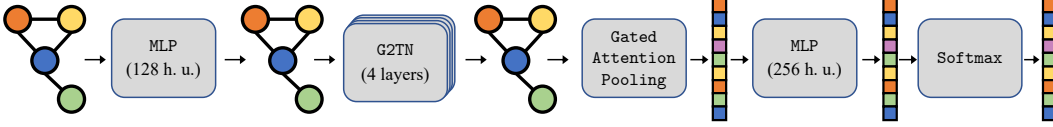


Figure 3: Visualization of the architecture used for NCI1 and NCI109 described in Section 5.

- the length of random walks;
- the number of low-rank functionals;
- the maximal tensor degree;
- the number of iterations (layers).

Note that the first three hyperparameters can also potentially vary across different iterations.

H Experiments

We have implemented the low-rank algorithm for our layers given in Theorem 3 using Tensorflow, Keras and Spektral [25]. Code is available at <https://github.com/tgcsaba/graph2tens>. All experiments were run on one of 3 computing clusters that were overall equipped with 5 NVIDIA Geforce 2080Ti, 2 Quadro GP100, 2 A100 GPUs. The largest GPU memory allocation any one of the experiments required was around ~ 5 GB. For the experiments ran using different random seeds, the seed was used to control the randomness in the (1) data splitting process, (2) parameter initialization, (3) optimization procedure. For an experiment with overall n_{runs} number of runs, the used seeds were $\{0, \dots, n_{\text{runs}} - 1\}$.

H.1 Model details.

The architecture described in Section 5 is visualized conceptually on Figure 3, where for G2T(A)N the G2TN stack is replaced by the attentional variation without changing any hyperparameters. Further to the discussion given in the main text, here we provide more details on the model implementation.

Initialization. In our G2TN and G2T(A)N layers, we linearly project tensor-valued features using linear functionals that use the same low-rank parametrization as the *recursive* variation in [62, App. D.2]. There, the authors also propose in Appendix D.5 an initialization heuristic for setting the variances of these component vectors, that we have employed with a centered uniform distribution.

Regularization. We also apply ℓ_2 regularization for both experiments in Sections 5 and H.3 for which we discuss the implementation here. Given a max. tensor degree $M \geq 1$ and a max. tensor rank $R \geq 1$ representing the layer width, there are overall RM rank-1 linear functionals of the form

$$\ell_m^r = u_{M-m+1}^r \otimes \dots \otimes u_M^r \quad \text{for } m = 1, \dots, M \text{ and } r = 1, \dots, R. \quad (28)$$

Since in practice these are represented using only the component vectors $u_i^r \in \mathbb{R}^d$, a naive application of ℓ_2 regularization would lead to computing the 2-norm of each component vector u_i^r to penalize the loss function. However, we found this approach to underperform compared to the following idea. Although our algorithm represents the functionals using a rank-1 decomposition and computes the projection of each tensor-valued node feature without explicitly building high-dimensional tensors, conceptually we still have tensors (ℓ_m^r) acting on tensors ($\Phi(i)$ for $i \in \mathcal{V}$), and hence the tensor norm, $\|\ell_m^r\|_2$, should be used for regularization. Fortunately, this can also be computed efficiently:

$$\|\ell_m^r\|_2 = \|u_{M-m+1}^r \otimes \dots \otimes u_M^r\|_2 = \|u_{M-m+1}^r\|_2 \cdots \|u_M^r\|_2.$$

Further, as is common, we replace the ℓ_2 norm with its squared value, sum over all functionals in (28), and multiply by the regularization parameter $\lambda > 0$ so that the final penalty is given by

$$\text{L2Penalty} = \lambda \sum_{r=1}^R \sum_{m=1}^M \|u_{M-m+1}^r\|_2^2 \cdots \|u_M^r\|_2^2.$$

H.2 Experiment Details.

Here we provide further details on the main experiment described in Section 5.

Hyperparameter Selection. The model architecture was primarily motivated by GraphTrans (small) from [70]. Specifically, the number and width of GNN layers, and the dropout rate was adopted as is. The ℓ_2 regularization strength was chosen equal to the weight decay rate. The other hyperparameters were tuned on a single split of NCI1. For the random walk length, we experimented with values 2, 5, 10. For the given GNN depth (4), 2 RW steps per layer was not enough to learn long-range interactions, while 10 significantly slowed down the convergence rate during training. For the maximal degree of tensors, we experimented with values from 2, 3, 4, and using values beyond 2 did not provide improvements. Intuitively, the tensor degree represents the order of nonlinear interactions that are learnable by the layer, e.g. a degree of 2 encodes pairwise interactions between node features in the neighbourhood, while a higher degree of M allows to encode interactions between certain M -tuples of nodes. We suggest that a degree of 2 is a good baseline setting, and that increasing the GNN depth instead allows to efficiently capture higher order interactions, while increasing the effective influence radius at the same time. For the pre- and postprocessing layers, we experimented with various depths and choosing more than 1 layer for each was counterproductive. The number of units was simply set to the GNN width (128) in the preprocessing layer, while for the postprocessing layer slightly increasing it was found to provide improvements (256), potentially to compensate for the large amount of information that is compressed in the pooling step.

Table 2: Accuracies computed over 5 seeds of G2T(A)N ablated by changing a single option.

Dataset	NoDiff	NoZeroStart	NoAlgOpt	NoJK	NoSkip	NoNorm	AvgPool
NCI1	80.1 \pm 0.7	79.5 \pm 1.8	81.6 \pm 1.6	82.1 \pm 1.8	81.8 \pm 0.9	81.6 \pm 1.5	82.4 \pm 0.9
NCI109	78.2 \pm 1.2	77.7 \pm 1.8	77.5 \pm 1.3	77.6 \pm 1.2	78.3 \pm 1.3	79.8 \pm 1.4	77.6 \pm 1.3

Ablation Studies. Further to using attention, we give a list of ideas for variations on our models in Appendix G, which can be summarized briefly as: (i) using increments of node attributes (Diff), (ii) prepending a zero point to sequences (ZeroStart), (iii) optimizing over the algebra embedding (AlgOpt), all of which are built into our main models. Further, the previous architectural choices aimed at incorporating several commonly used tricks for training GNNs. We investigate the effect of the previous variations and ablate the architectural “tricks” by measuring the performance change resulting from ceteris paribus removing it. Table 2 shows the result for G2T(A)N. To summarize the main observations, the model is robust to all the architectural changes, removing the layer norm even improves on NCI109. Importantly, replacing the attention pooling with mean pooling does not significantly affect the performance, but actually slightly improves on NCI1. Regarding variations, AlgOpt slightly improves on both datasets, while removing Diff and ZeroStart significantly degrades the accuracy on NCI1. The latter means that translations of node attributes are important, not just their distances.

We give the analogous ablation study for the G2TN model in Table 3, and compare the derived conclusions between the attentional and attention-free versions. First, we discuss the layer variations. Similarly to G2T(A)N, NoDiff slightly decreases the accuracy. However the conclusions regarding NoZeroStart and NoAlgOpt are different. In this case, removing ZeroStart actually improves the performance, while on G2T(A)N the opposite was true. An interpretation of this phenomenon is that only the attention mapping that is used to learn the random walks required information about translations of the layer’s features, and not the tensor-features themselves. Another difference is that NoAlgOpt degrades the accuracy more significantly for G2TN. A possible explanation is that since G2T(A)N layers are more flexible thanks to their use of attention, they rely less on being able to learn the algebraic lift, while as G2TN layers are more rigid in their random walk definition, and benefit more from the added flexibility of AlgOpt. Additionally, it seems that G2TN is more sensitive to the various architectural options, and removing any of them, i.e. NoJK, NoSkip, NoNorm or AvgPool, degrades the accuracy by 1% or more on at least one dataset. Intuitively, it seems that overall the G2T(A)N model is more robust to the various architectural “tricks”, and more adaptable due to its ability to learn the the random walk.

Table 3: Accuracies computed over 5 seeds of G2TN ablated by changing a single option.

Dataset	NoDiff	NoZeroStart	NoAlgOpt	NoJK	NoSkip	NoNorm	AvgPool
NCII	79.1 ± 1.5	80.7 ± 0.6	78.9 ± 1.3	79.4 ± 1.9	81.0 ± 1.3	79.5 ± 1.1	80.3 ± 1.7
NCII09	77.8 ± 1.2	79.5 ± 1.5	76.5 ± 1.7	78.1 ± 1.9	77.0 ± 1.6	77.4 ± 2.7	77.6 ± 1.8

Table 4: Accuracies of our models on the citation datasets computed over 100 seeds compared with the 4 consistently best performing baselines from [57].

Dataset	GCN	GAT	MoNet	GraphSage (mean)	G2TN (ours)	G2T(A)N (ours)
Cora	81.5 ± 1.3	81.8 ± 1.3	81.3 ± 1.3	79.2 ± 7.7	82.6 ± 1.0	82.0 ± 1.1
Citeseer	71.9 ± 1.9	71.1 ± 1.9	71.2 ± 2.0	71.6 ± 1.9	69.4 ± 1.0	68.2 ± 1.3
Pubmed	77.8 ± 2.9	78.7 ± 2.3	78.6 ± 2.3	77.4 ± 2.2	78.8 ± 1.9	78.0 ± 1.9

H.3 Further Experiments

Citation datasets. Additionally to transductive learning on the biological datasets, we have carried out inductive learning tasks on some of the common citation datasets, i.e. Cora, Citeseer [56] and Pubmed [47]. We follow [57] in carrying out the experiment, and use the largest connected component for each dataset with 20 training examples per class, 30 validation examples per class, and the rest of the data used for testing. The hyperparameters of our models and optimization procedure were based on the settings of the GAT model in [57, Table 4], which we have slightly fine-tuned on Cora and used for the other datasets. In particular, a single layer of G2TN or G2T(A)N is used with 64 functionals, max. tensor degree 2 and random walk length 5. The dropout rate was tuned to 0.9, while the attentional dropout was set to 0.3 in G2T(A)N. Optimization is carried out with Adam [32], a fixed learning rate of 0.01 and ℓ_2 regularization strength 0.01. Training is stopped once the validation loss does not improve for 50 epochs, and restored to the best checkpoint. For both of our models, NoDiff is used that we found to improve on the results as opposed to using increments of node features. The dropout rate had to be tuned as high as 0.9, which suggests very strong overfitting, hence the additional complexity of AlgOpt was also contrabeneficial, and NoAlgOpt was used. As such, each model employs a single hidden layer, which is followed by layer normalization that was found to perform slightly better than other normalizations, e.g. graph-level normalization.

The results of our models over 100 runs are reported in Table 4 compared with the 4 consistently best performing baselines on these datasets from [57], i.e. GCN [33], GAT [65], MoNet [45], and GraphSage [29] with a mean aggregator. Firstly on Cora, both G2TN and G2T(A)N outperform the baselines with a more significant improvement for G2TN. For CiteSeer, our models are left somewhat behind compared to the baselines in terms of accuracy. Finally, they are again competitive on Pubmed, where G2TN takes the top score with a very slight lead. Two consistent observations are: (1) G2TN and G2T(A)N have a lower variance than all baselines, (2) G2TN consistently outperforms G2T(A)N. The latter may be attributed to the observation that due to the severe overfitting on these datasets, the additional complexity of the attention mechanism in G2T(A)N is unhelpful for generalization.

K-hop Sanitized Splits. Recent work [50] has demonstrated that it is possible to make the previously considered citation datasets more suitable for testing the ability of a model to learn long-range information by dropping node labels in a structured way. Concretely, they use a label resampling strategy to guarantee that if a node is selected for a data split, none of its k -th degree neighbours are included in any splits, i.e. training, validation nor testing, allowing to reduce the effect of short-range “label imprinting”. In practice, we select a maximal independent set from the graph with respect to the k -th power of the adjacency matrix with self-loops, and repeat the previous experiment with the same data splitting method, model choice and training settings as before. In this case, the experiment seed is also used to control the random maximal independent set that is selected.

The results of our models trained on the citation datasets sanitized this way are available in Table 5 computed over 100 seeds for $k = 1, 2$. As baseline results, we compare against the 5 best performing models from [50], where various GNN depths were also considered for each model, and we use the *best* reported result for each of the baselines. Overall, we can observe that all baseline models exhibit a very sharp drop in performance as k is increased, while for G2TN and G2T(A)N, the performance

Table 5: Accuracies of our models on k -hop sanitized citation datasets computed over 100 seeds compared with the 5 consistently best performing models from [50].

k	Dataset	GCN	GAT	g-U-net	HGNet-EP	HGNet-L	G2TN (ours)	G2T(A)N (ours)
1	Cora	76.7	78.5	78.1	77.2	77.1	81.8 ± 1.3	80.9 ± 1.4
	Citeseer	64.2	66.4	63.0	64.3	64.10	68.1 ± 1.3	66.6 ± 1.4
	Pubmed	75.8	75.9	75.8	77.0	76.3	78.7 ± 1.9	77.5 ± 2.0
2	Cora	72.0	73.4	74.4	74.0	75.4	79.4 ± 2.8	78.3 ± 2.9
	Citeseer	58.3	59.4	57.3	57.8	59.9	63.6 ± 1.8	62.2 ± 1.9
	Pubmed	72.1	73.1	72.4	72.9	75.1	77.1 ± 1.7	76.3 ± 1.6

Table 6: Computation time of one forward pass of one G2TN layer in milliseconds on a Nvidia GeForce 2080TI GPU for a graph with varying nodes (N) and edges (E), while node feature dimension is fixed at 128.

	$E = 5000$	$E = 10000$	$E = 20000$	$E = 40000$	$E = 80000$
$N = 500$	37	61	98	151	293
$N = 1000$	37	62	104	152	290
$N = 2000$	39	63	89	160	293
$N = 4000$	39	63	109	160	294
$N = 8000$	42	65	108	156	295

decrease is not nearly as pronounced. For both $k = 1$ and $k = 2$, our models perform better than the baselines on all datasets. This is explained by the fact that due to using random walks length of 5, the models can efficiently pick up on information outside of the sanitized neighbourhoods. As previously, G2TN performs better than G2T(A)N as the additional flexibility of the attention layer does not lead to improved generalization performance when severe overfitting is present. This experiment demonstrates that the proposed models efficiently pick up on long-range information within larger neighbourhoods, and it suggests that on inductive learning tasks they should be more robust to sparse labeling rates compared to common short-range GNN models.

H.4 Computation Time

Table 6 shows the computation time of one forward pass of a G2TN layer with various graphs. This empirically demonstrates that our layer does not depend on the number of nodes, and only depends linearly on the the number of edges as expected from the theoretical complexity.

Table 7 shows the computation time of one forward pass of a G2TN layer with a fixed graph and various model parameters. Empirically, our layer scales roughly linearly in both the walk length k and the maximum tensor level M . Linear complexity in walk length is expected, but the linearity in the maximum tensor level M is due to parallelization of the algorithm presented in Theorem 3. This is discussed in the pseudocode section of Appendix F.

H.5 Parameter Count

In both the NCI datasets, the G2TN and G2T(A)N have 505k and 519k trainable parameters respectively. Thus, our models are comparable in size to the GraphTrans, which is reported to have 500k trainable parameters [70].

H.6 Summary of Dataset Parameters

Here, we provide a summary of the dataset statistics and the model parameters used in the experiments. Here, M denotes the maximum tensor degree, k denotes the walk length, and R denotes the maximum rank of the low rank functions. Concretely, R is the number of node features used within the network, and is analogous to the width of a neural network. For the computation of the diameter of a disconnected graph, we take the largest diameter of any connected component.

Table 7: Computation time of one forward pass of one G2TN layer in milliseconds on a Nvidia GeForce 2080TI GPU for a fixed graph with $N = 2000$ nodes and $E = 10000$ edges. The walk length (k) and the maximum level M of the model are varied.

	$k = 2$	$k = 4$	$k = 6$	$k = 8$	$k = 10$
$M = 1$	14	26	45	59	70
$M = 2$	20	40	72	81	101
$M = 3$	26	55	95	115	163
$M = 4$	34	73	101	150	188
$M = 5$	42	92	139	187	237

Table 8: Dataset statistics and G2TN model parameters used.

	graphs	total nodes	total edges	avg. diam.	classes	layers	M	k	R
NCI1	4110	122k	132k	13.3	2	4	2	5	128
NCI109	4127	122k	132k	13.1	2	4	2	5	128
CORA	1	2485	5069	19	7	4	2	5	64
Citeseer	1	2110	3668	28	6	4	2	5	64
Pubmed	1	19k	44k	18	3	4	2	5	64