

MULTI-CHANNEL GRAPH CONVOLUTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Defining the convolution on graphs has led to much progress in graph machine learning, particularly through approximations based on polynomials and, ultimately, message-passing neural networks (MPNNs). However, this convolution is defined for single-channel graph signals, i.e., a single feature is given at each node, and a single new feature is assigned to each node. As multiple initial node features are provided for many challenging tasks and convolutions are generally defined for these multi-channel signals, we introduce multi-channel graph convolutions (MGCs) by obtaining their form using the graph Fourier transform. MGCs highlight the critical importance of utilizing multiple edge relations to amplify different signals for each feature channel. We further introduce localized multi-channel MPNNs and the multi-channel graph isomorphism network (MC-GINs), with which we can provably obtain linear mappings that are injective on multisets. Our experiments confirm the greatly improved capabilities of MGCs and MC-GINs.

1 INTRODUCTION

Many challenging tasks and applications are based on graph-structured data, e.g., property prediction for molecules (Hu et al., 2021), fraud detection in transactions (Weber et al., 2019), and recommendations (Monti et al., 2019). Developing expressive and well-performing methods to learn from such data is thus an important challenge. As one such method, neural networks for graph-structured data are based on defining a convolution on a graph. Spectral graph convolutions emerged by computing the convolution exactly in the graph Fourier domain based on the convolution theorem (Hammond et al., 2011; Bruna et al., 2014). Due to its high computational cost, various approximations based on polynomials of the graph Laplacian emerged (Levie et al., 2019; He et al., 2021; Koke & Cremers, 2024), e.g., using Chebyshev polynomials (Defferrard et al., 2016). This further led to the graph convolutional network (GCN) (Kipf & Welling, 2017) as a localized approximation. Most currently used message-passing neural networks (MPNNs) are derivations and improvements of this spectral graph convolution and the GCN. This includes applying different aggregation functions like the mean (Hamilton et al., 2017), the sum (Xu et al., 2019), utilizing attention coefficients (Velickovic et al., 2018; Brody et al., 2022) or negative edge weights (Yan et al., 2022). More complex methods like gating mechanisms (Li et al., 2016; Rusch et al., 2023), positional encodings (Kreuzer et al., 2021; Rampásek et al., 2022; Huang et al., 2024), and normalization layers (Zhao & Akoglu, 2020) are typically combined with these approximated graph convolutions.

However, this convolution is defined for single-channel signals, i.e., each node has a single feature and a single feature per node is obtained. In most applications, each node has multiple initial features assigned to it, e.g., a text embedding for documents and atom features for molecular data. Similarly, the goal of these models is often to find rich node embeddings capturing both structural properties and feature interactions, which are typically designed to have multiple feature channels. The currently used definition of graph convolutions and, consequently, most MPNNs are ill-defined for this task. As MPNNs amplify the same signal for each channel, issues like representational rank collapse (Roth & Liebig, 2023) and over-smoothing (Oono & Suzuki, 2020) emerged.

The convolution that is mapping multi-channel signals to multi-channel signals is generally defined differently. We introduce these multi-channel convolutions to graphs by obtaining their form based on the convolution theorem and the graph Fourier transform. This multi-channel graph convolution (MGC) highlights the importance of utilizing multiple edge relations. Each edge relation corresponds to a different signal, and the feature transformation describes the amplification or damping

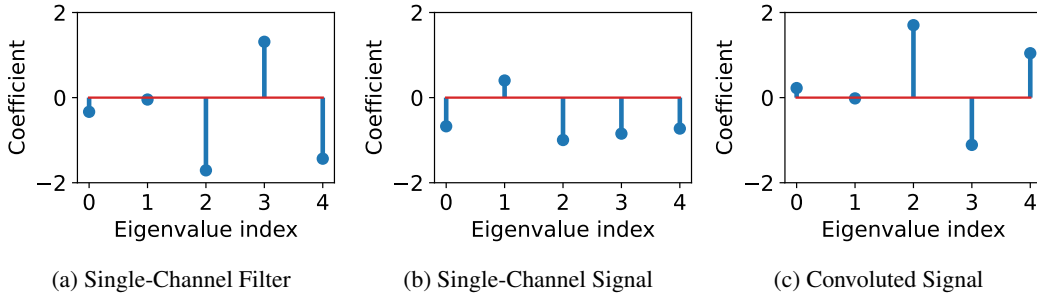


Figure 1: Single-Channel signal and filter in the Fourier domain. The eigenvalue index corresponds to the eigenvalues (frequencies) of the graph Laplacian L . The convoluted single-channel signal in the Fourier domain is the element-wise scalar product of the filter and signal.

of this signal for each output feature channel. While the exact formulation is computationally prohibitive for large graphs, knowing the exact form of the graph convolution for multiple channels will allow for various approximations. We define required properties for localized approximations that similarly operate on multiple sparse edge relations. We further introduce the multi-channel graph isomorphism network (MC-GIN) that we prove to have the same expressivity as the graph isomorphism network (Xu et al., 2019) while applying a linear transformation to the data. Our experiments confirm the advantages of the MCGC and MC-GIN. We summarize our main contributions as follows:

- Based on the general definition of the convolution for multi-channel signals, we obtain the multi-channel graph convolution (MCGCs) using the convolution theorem (Section 3).
- To construct multi-channel MPNNs (MC-MPNNs), we introduce localized MCGCs, for which we require multiple edge relations, with each edge relation amplifying a different signal in the data (Section 3.1).
- We prove that with localized MCGCs, we can construct linear mappings that are injective on multisets, which we call multi-channel graph isomorphism network (MC-GIN) due to its expressivity being equivalent to the GIN (Xu et al., 2019) (Section 3.2).

2 PRELIMINARIES

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a connected and undirected graph consisting of a set of n nodes \mathcal{V} and a set of edges \mathcal{E} . A graph signal is defined as a function $x: \mathcal{V} \rightarrow \mathbb{R}^d$ that assigns a vector of real values to each node. For notational simplicity, we stack all node signals into a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ based on some node ordering. For $d = 1$, we refer to this as a single-channel signal, while for $d > 1$, we call it a multi-channel signal. These can either be initial features or expressive and informative node embeddings obtained by a suitable method. Let $\mathbf{A} \in \{0, 1\}^{n \times n}$ with $A_{i,j} = 1$ if $(i, j) \in \mathcal{E}$ and 0 otherwise be the adjacency matrix corresponding to the same node ordering as \mathbf{X} . The diagonal degree matrix is $\mathbf{D} \in \mathbb{N}^{n \times n}$. The symmetrically normalized adjacency matrix is given by $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ and the graph Laplacian by $\mathbf{L} = \mathbf{I}_n - \tilde{\mathbf{A}}$. Its eigendecomposition is $\mathbf{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ where $\mathbf{\Lambda} \in \mathbb{R}^{n \times n}$ is a diagonal matrix containing its eigenvalues, and $\mathbf{U} \in \mathbb{R}^{n \times n}$ is an orthonormal matrix containing the corresponding eigenvectors as columns. In the graph domain, the Fourier base is given by the eigenvectors \mathbf{U}^T of the graph Laplacian. Thus, the Fourier transformation $F = \mathbf{U}^T$ is performed by projecting a graph signal onto the eigenvectors, and its inverse transformation is given by $F^{-1} = \mathbf{U}$.

2.1 SINGLE-CHANNEL GRAPH CONVOLUTIONS

The convolution is an operation that combines two functions and produces a new function. Given a graph with n nodes, the convolution is typically defined as the function $x' = w * x$, where $x \in \mathbb{R}^n$ is a single-channel graph signal, $w \in \mathbb{R}^n$ is a corresponding filter, and it produces a convoluted single-channel graph signal $w' \in \mathbb{R}^n$. As this convolution involves single-channel signals, we will further

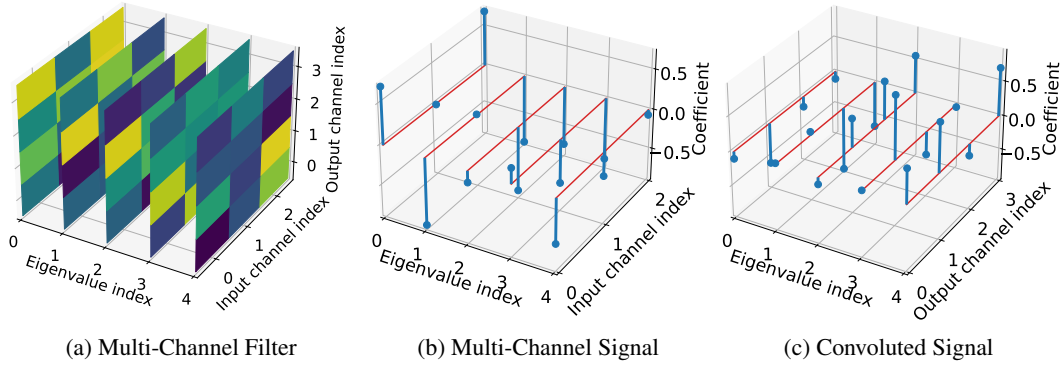


Figure 2: Multi-channel signal and filter in the Fourier domain. The eigenvalue index corresponds to the eigenvalues (frequencies) of the graph Laplacian L . The element-wise product for each eigenvalue index is a matrix-vector product.

refer to it as the single-channel graph convolution (SCGC). The convolution theorem (O’Neil, 1963) states that the Fourier transform of a convolution

$$F(\mathbf{w} * \mathbf{x}) = F(\mathbf{w}) \odot F(\mathbf{x})$$

is equal to the element-wise multiplication of the filter and signal in the Fourier domain. The graph convolution (Hammond et al., 2011) can be expressed as

$$\begin{aligned} \mathbf{w} * \mathbf{x} &= F^{-1}(F(\mathbf{w}) \odot F(\mathbf{x})) \\ &= U(U^T \mathbf{w} \odot U^T \mathbf{x}) \end{aligned}$$

by utilizing the graph Fourier transform $F = U^T$ and its inverse $F^{-1} = U$ based on the eigenvectors U of the graph Laplacian L . We visualize this process in Figure 1. Substituting $U^T \mathbf{w}$ with its diagonalized matrix $\mathbf{W}^* = \text{diag}(U^T \mathbf{w})$ and the Hadamard product with a matrix multiplication leads to the equivalent form

$$U(U^T \mathbf{w} \odot U^T \mathbf{x}) = U\mathbf{W}^*U^T \mathbf{x}.$$

Bruna et al. (2014) proposed to learn the filter \mathbf{W}^* directly in the Fourier domain, which is also known as the spectral graph convolution.

2.2 APPROXIMATIONS USING POLYNOMIALS

As computing the eigendecomposition and performing dense matrix multiplications is computationally expensive, SCGCs do not scale well to large graphs. Defferrard et al. (2016) proposed to instead approximate and parameterize graph convolutions using a learnable function $g(\Lambda)$ on the eigenvalues Λ of L , i.e., by approximating the diagonal matrix $\mathbf{W}^* \approx g(\Lambda)$. One commonly employed approximation

$$\begin{aligned} U\mathbf{W}^*U^T &\approx U \sum_{p=1}^n w_p U T_p(\tilde{\Lambda}) U^T \\ &= \sum_{p=0}^k w_p T_p(U\Lambda^*U^T) \\ &= \sum_{p=0}^k w_p T_p(\tilde{L}) \end{aligned}$$

is based on Chebyshev polynomials (Chebyshev, 1853) up to k -th order of a rescaled matrix of eigenvectors $\tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - I_n$ and corresponding $\tilde{L} = \frac{1}{\lambda_{\max}} L - I_n$ (Hammond et al., 2011). The polynomials are defined as $T_0(\Lambda) = I$, $T_1(\Lambda) = \Lambda$, and $T_k(\Lambda) = 2\Lambda T_{k-1}(\Lambda) - T_{k-2}(\Lambda)$. This

truncated polynomial expansion is k -localized in the graph as entries $\left[\sum_{p=0}^k w_p T_p(\tilde{\mathbf{L}})\right]_{i,j}$ are zero when the shortest path between nodes i and j is larger than k . Further approximations include Cayley polynomials (Levie et al., 2019), Bernstein polynomials (He et al., 2021). Koke & Cremers (2024) propose an approximation for directed graphs using Faber polynomials. GPR-GNN proposes to directly learn the filter weights (Chien et al., 2021). All these methods are approximations of the SCGC.

2.3 MESSAGE-PASSING NEURAL NETWORKS

For further computational efficiency and empirical success, the graph convolutional network (GCN) (Kipf & Welling, 2017) was derived as a 1-localized approximation of Chebyshev polynomials. It uses the approximation

$$\mathbf{w} * \mathbf{x} \approx w_0(\mathbf{I}_n + \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}) \mathbf{x}.$$

by setting $k = 1$, $\lambda_{\max} \approx 2$, and $w_0 = -w_1$. They further substitute $\mathbf{I}_n + \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ by $\tilde{\mathbf{D}}^{-1/2} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-1/2}$ using $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_n$ and $\tilde{\mathbf{D}} = \mathbf{D} + \mathbf{I}_n$.

However, as the SCGC is not defined for multi-channel signals, the SCGC and its approximations are not directly applicable to multi-channel signals. To apply the GCN to a multi-channel signal $\mathbf{X} \in \mathbb{R}^{n \times d}$ or to obtain a multi-channel signal $\mathbf{X}' \in \mathbb{R}^{n \times d'}$ as the output, they propose to replace the scalar $w \in \mathbb{R}$ with a matrix $\mathbf{W} \in \mathbb{R}^{d \times d'}$ (Kipf & Welling, 2017). This leads to their update function

$$\mathbf{X}' = \hat{\mathbf{A}} \mathbf{X} \mathbf{W}, \tag{1}$$

where $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-1/2} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-1/2}$.

Recent research on over-smoothing and rank collapse has proven that graph convolutions of the form in Eq. 1 do not allow the amplification of different signals across channels (Giovanni et al., 2023; Roth & Liebig, 2023; Roth, 2024). The signals to be amplified are solely determined by the spectrum of $\hat{\mathbf{A}}$, which was shown to hold for any choice of $\hat{\mathbf{A}}$ and \mathbf{W} . While this may be desired if the same signal should be amplified for all feature columns, each feature channel should typically be composed of a different mix of signals. Thus, the GCN inherited this issue from SCGCs.

Most established message-passing neural networks (MPNNs) are adaptations of the GCN and similarly suffer from being based on single-channel convolutions. Examples include utilizing the mean (Hamilton et al., 2017) or sum (Xu et al., 2019) aggregation, using the attention mechanism (Velickovic et al., 2018; Brody et al., 2022), allowing for negative edge weights (Yan et al., 2022) or building on top of these graph convolutions with normalization layers (Zhao & Akoglu, 2020), gating mechanisms (Li et al., 2016; Rusch et al., 2023), or positional encodings (Kreuzer et al., 2021; Rampásek et al., 2022; Huang et al., 2024). We will thus consider how to directly define the graph convolution for multi-channel signals, of which all approximations would benefit.

3 MULTI-CHANNEL GRAPH CONVOLUTIONS

We now consider a multi-channel graph signal $\mathbf{X} \in \mathbb{R}^{n \times d}$ with d channels for each node. By applying a graph convolution, we want to obtain a convoluted multi-channel graph signal with different channel mixing for any input and output combination. In the last section, we have shown that the currently used SCGC is not defined for this task. In general, the convolution is defined for multi-channel signals \mathbf{X} . To obtain a convoluted signal with c channels, the element-wise product needs to map vectors with d channels to vectors with c channels. This requires elements of the filter to be matrices $\mathbf{W}_i \in \mathbb{R}^{c \times d}$ and the full filter to be $\mathbf{W} \in \mathbb{R}^{n \times c \times d}$. This is commonly referred to as the multi-channel convolution Burg (1964); Inouye & Sato (1999). It is commonly used in signal processing (Burg, 1964; Inouye & Sato, 1999; Sainath et al., 2017), e.g., for multi-channel Wiener filtering (Brandstein & Ward, 2001) and finite impulse responses (Seltzer et al., 2004). Convolutional neural networks (CNNs) (LeCun et al., 1989) are similarly based on the multi-channel convolution (Zhang et al., 2021).

Equivalently to the single-channel convolution, we can compute the multi-channel graph convolution in the Fourier domain using the convolutional theorem. The exact form of the multi-channel

convolution on graphs is a direct consequence of the general multi-channel convolution and the graph Fourier transform:

Theorem 3.1. (*Multi-Channel Graph Convolution*) Let $\mathbf{L} \in \mathbb{R}^{n \times n}$ be a graph Laplacian of an undirected graph and the graph Fourier transform be given by its matrix of eigenvectors $\mathbf{U}^T \in \mathbb{R}^{n \times n}$. Further, let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a multi-channel graph signal, and $\mathbf{W} \in \mathbb{R}^{n \times c \times d}$ be a corresponding filter matrix. Then,

$$\mathbf{W} * \mathbf{X} = \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T \mathbf{X} \hat{\mathbf{W}}_i \quad (2)$$

where $\hat{\mathbf{W}}_k = [\mathbf{U}^T \mathbf{W}]_k^T$ is the graph Fourier transformed filter, \mathbf{u}_k is the k -th eigenvector of \mathbf{L} .

Proof. We use the vectorized signal $\hat{\mathbf{x}} = \text{vec}(\mathbf{X}) \in \mathbb{R}^{n \cdot c}$ by stacking its columns. The graph Fourier transform on matrices and tensors is applied along the node dimension, i.e., independently on each channel. For matrix \mathbf{X} this results in

$$F(\mathbf{X}) = \mathbf{U}^T \mathbf{X} \in \mathbb{R}^{n \times d}$$

and for tensor \mathbf{W} in

$$\hat{\mathbf{W}} = F(\mathbf{W}) = \mathbf{U}^T \times_1 \mathbf{W} \in \mathbb{R}^{n \times c \times d}$$

where \times_1 is the 1-mode tensor matrix product (Kolda & Bader, 2009) that performs the desired broadcasted matrix multiplication. With this, we state the multi-channel graph convolution in the Fourier domain as

$$\mathbf{W} * \mathbf{X} = \mathbf{U}(\mathbf{U}^T \times_1 \mathbf{W} \odot \mathbf{U}^T \mathbf{X}) = \mathbf{U}(\hat{\mathbf{W}} \odot \mathbf{U}^T \mathbf{X}).$$

This product is visualized in Figure 2. Similarly to the single-channel, we simplify this expression using matrix multiplications. Equivalently to the single-channel case, this can be achieved by diagonalizing $\hat{\mathbf{W}}$ into a block matrix of diagonal blocks

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{W}}_{1,0,0} & 0 & 0 & \hat{\mathbf{W}}_{1,0,d} & 0 & 0 \\ 0 & \ddots & 0 & \dots & 0 & \ddots & 0 \\ 0 & 0 & \hat{\mathbf{W}}_{n,0,0} & 0 & 0 & \hat{\mathbf{W}}_{n,0,d} \\ \vdots & & & \ddots & & \vdots & \\ \hat{\mathbf{W}}_{1,c,0} & 0 & 0 & \hat{\mathbf{W}}_{1,c,d} & 0 & 0 \\ 0 & \ddots & 0 & \dots & 0 & \ddots & 0 \\ 0 & 0 & \hat{\mathbf{W}}_{n,c,0} & 0 & 0 & \hat{\mathbf{W}}_{n,c,d} \end{bmatrix} \in \mathbb{R}^{nc \times nd}$$

where each $\hat{\mathbf{W}}_{k,i,j} \in \mathbb{R}$ corresponds to position i, j in $\hat{\mathbf{W}}_k \in \mathbb{R}^{c \times d}$. This simplifies the equivalent vectorized form into

$$\text{vec}(\hat{\mathbf{W}} \odot \mathbf{U}^T \mathbf{X}) = \mathbf{D} \text{vec}(\mathbf{U}^T \mathbf{X}) = \mathbf{D} (\mathbf{I}_d \otimes \mathbf{U}^T) \text{vec}(\mathbf{X})$$

by utilizing the Kronecker product \otimes . The matrix \mathbf{D} can further be decomposed into a sum of Kronecker products $\mathbf{D} = \sum_{i=1}^n \hat{\mathbf{W}}_i \otimes \mathbf{I}_n^{(i)}$, where for each $\mathbf{I}_n^{(i)} \in \mathbb{R}^{n \times n}$ all entries are zero, apart from position i, i which is one. This lets us state the full vectorized multi-channel graph convolution as

$$\begin{aligned} \text{vec}(\mathbf{W} * \mathbf{X}) &= (\mathbf{I}_n \otimes \mathbf{U}) \left(\sum_{i=1}^n \hat{\mathbf{W}}_i \otimes \mathbf{I}_n^{(i)} \right) (\mathbf{I}_n \otimes \mathbf{U}^T) \text{vec}(\mathbf{X}) \\ &= \left(\sum_{i=1}^n \hat{\mathbf{W}}_i \otimes \mathbf{u}_i \mathbf{u}_i^T \right) \text{vec}(\mathbf{X}) \end{aligned}$$

by using the fact that $\mathbf{U} \mathbf{I}_n^{(i)} \mathbf{U}^T = \mathbf{u}_i \mathbf{u}_i^T$. Inverting the vec operation allows us to avoid the Kronecker product and state the exact multi-channel graph convolution as

$$\mathbf{W} * \mathbf{X} = \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^T \mathbf{X} \hat{\mathbf{W}}_i^T$$

This concludes the proof. \square

We emphasize that this form is a mathematical fact and not a definition made by us. We interpret the multi-channel graph convolutions as follows: Each term $\mathbf{u}_i \mathbf{u}_i^T \mathbf{X} \mathbf{W}_i$ corresponds to one Fourier basis vector. The corresponding parameter matrix \mathbf{W}_i specifies how much this signal coming from each input channel should be amplified or damped for each output channel. Utilizing n terms allows the multi-channel graph convolution to amplify and filter specific input signals for each output channel.

Now that we have obtained the graph convolution for multi-channel input or output signals, we could directly apply it to suitable tasks. However, we expect similar issues for the multi-channel graph convolution as we have seen for the single-channel graph convolution. Issues with directly computing the single-channel graph convolution include high computational costs due to dense matrix multiplication and computing the eigendecomposition. It is also inherently transductive, as the eigenvectors change across graphs, and thus, a learned filter is not applicable. It was also observed that the single-channel graph convolution quickly overfits given training data. Many directions have been proposed to approximate the single-channel convolution, as outlined in Section 2.1. The most prominent approximation and derivation from single-channel graph convolutions are MPNNs. While a similar multitude of novel approximations can be constructed for the multi-channel graph convolution, we want to outline how MPNNs based on MCGC should be designed.

3.1 MULTI-CHANNEL MPNNs

Each outer product $\mathbf{u}_i \mathbf{u}_i^T$ can be seen as one edge relation type with full connectivity and edge weights $u_{i,p} u_{i,q}$ between nodes p and q . Multi-channel MPNNs should be localized and still utilize multiple edge relations $\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(k)} \in \mathbb{R}^{n \times n}$, with each relation amplifying a different signal. While a sparse graph will have more than one non-zero eigenvalue, its dominating eigenvector can be controlled by its edge weights. The corresponding transformations $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)} \in \mathbb{R}^{d \times d'}$ can then similarly determine the amplification or damping of these signals for each channel. Splitting a graph into a multi-relational graph was recently proposed to avoid representational rank collapse (Roth et al., 2024). We define localized multi-channel graph convolutions as follows:

Definition 3.2. (Localized Multi-Channel Graph Convolution (l-MCGC)) Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be an adjacency matrix. A function

$$\phi(\mathbf{X}) = \mathbf{X}' = \sum_{m=1}^k \hat{\mathbf{A}}^{(m)} \mathbf{X} \mathbf{W}^{(m)} \quad (3)$$

with $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(k)} \in \mathbb{R}^{n \times n}$, and $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)} \in \mathbb{R}^{d \times c}$ is called a localized multi-channel graph convolution (l-MCGC) if the following two conditions are satisfied:

- the dominant eigenvectors of $\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(k)}$ are linearly independent and
- $\forall i, j \in \{1, \dots, n\}, l \in \{1, \dots, k\} : \mathbf{A}_{ij} = 0 \implies \hat{\mathbf{A}}_{ij}^{(l)} = 0$.

These l-MCGCs can be equivalently expressed as a node-based update function

$$\begin{aligned} \mathbf{x}'_i = \phi(\mathbf{X})_i &= \sum_{m=1}^k \sum_{j \in \mathbb{N}_i} \hat{a}_{i,j}^{(m)} \mathbf{x}_j \mathbf{W}^{(m)} \\ &= \sum_{j \in \mathbb{N}_i} \mathbf{x}_j \mathbf{W}^{(i,j)} \end{aligned}$$

for node i using its set of neighbors \mathbb{N}_i , edge weights $a_{i,j}^{(m)} = \hat{a}_{i,j}^{(m)}$. For each message, a linear combination $\mathbf{W}^{(i,j)} = \sum_{m=1}^k \hat{a}_{i,j}^{(m)} \mathbf{W}^{(m)}$ of feature transformations $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)}$ is applied to the corresponding node state. This form also confirms that l-MCGCs remain permutation equivariant when requiring all $\hat{a}_{i,j}^{(m)}$ to be obtained using a permutation equivariant method. There are many ways to obtain a set of graphs $\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(k)}$ with distinct eigenvectors. As graph attention networks (Velickovic et al., 2018) with multiple attention heads can be expressed in this form as well, we will similarly refer to each term as a head. However, each $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(k)}$ is a row-stochastic

matrix because of the softmax activation. As all row-stochastic matrices have the constant vector as dominant eigenvector (Asmussen, 2003), all heads amplify the same signal. We now identify a connection between l-MCGCs and their expressivity with respect to the Weisfeiler-Leman test (Weisfeiler & Lehman, 1968), which leads to one possible instantiation of l-MCGCs.

3.2 EXPRESSIVITY

Our analysis has shown the advantages of multi-channel graph convolutions for amplifying different signals across channels. However, much attention has been paid to the structural expressivity of MPNNs, which studies the ability of MPNNs to distinguish non-isomorphic graphs. It was determined that MPNNs are upper-bounded in expressivity by the Weisfeiler-Leman test (Xu et al., 2019; Morris et al., 2019). To achieve this maximal expressivity, the graph isomorphism network (GIN) Xu et al. (2019) is typically employed and applies a complex non-linear feature transformation. We show that we can achieve the same expressivity by obtaining a linear l-MCGC:

Proposition 3.3. (*Injectivity on multisets.*) *Let $\mathbb{X} = \{\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\}$ be a countable multiset with $\mathbf{x}_i \in \mathbb{R}^{1 \times d}$ for $d \in \mathbb{N}$. Then, there exists a function $f: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^k$ so that*

$$\mathbf{x}'_i = \sum_{l=1}^k \sum_{\mathbf{x}_j \in \mathbb{X}_i} f(\mathbf{x}_i, \mathbf{x}_j)_l \mathbf{x}_j \mathbf{W}^{(l)}$$

is injective for all $\mathbf{x}_i \in \mathbb{X}$ and $\mathbb{X}_i \subset \mathbb{X}$ of bounded size, all $k > 1$, and a.e. $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)} \in \mathbb{R}^{d \times d'}$ with $d' \in \{1, \dots, n\}$.

Proof. Given any $\mathbf{x}_j \in \mathbb{X}$, two linear transformations $\mathbf{W}^{(m)}, \dots, \mathbf{W}^{(n)}$ map \mathbf{x}_j to pairwise linearly independent vectors $\mathbf{y}_j^{(m)} = \mathbf{x}_j \mathbf{W}^{(m)}$ and $\mathbf{y}_j^{(n)} = \mathbf{x}_j \mathbf{W}^{(n)}$ for a.e. $\mathbf{W}^{(m)}, \dots, \mathbf{W}^{(n)}$ with respect to the Lebesgue measure.

To prove the existence of the desired f , we follow the proof of Lemma 5 in Xu et al. (2019) (injectivity of GIN): Because \mathbb{X} is countable, there exist injective mappings $Z_1: \mathbb{X} \rightarrow \mathbb{N}$, $Z_2: \mathbb{X} \rightarrow \mathbb{N}$ mapping elements $\mathbf{x}_m \in \mathbb{X}$ to natural numbers. As each \mathbb{X}_n is of bounded size, there exists a number $N \in \mathbb{N}$ so that $|\mathbb{X}| < N$ for all X . An example of such f is $f(\mathbf{x}_m, \mathbf{x}_n) = [N^{-Z_1(\mathbf{x}_m)} \quad 1 \quad \dots \quad 1 \quad N^{-Z_2(\mathbf{x}_n)}]$. The first and last output values of f can be viewed as continuous one-hot vectors of \mathbf{x}_m and \mathbf{x}_n , respectively. Because \mathbb{X} is countable, Z_1 and Z_2 can further be chosen to output linearly independent vectors for all pairs of distinct inputs.

Thus, for all $\mathbf{x}_i \in \mathbb{X}$ and $\mathbf{x}_j \in \mathbb{X}$, we have a different linear combination of pairwise linear independent vectors and pairwise linearly independent linear combinations, and thus $\mathbf{s}_{i,j} = \sum_{l=1}^k f(\mathbf{x}_i, \mathbf{x}_j)_l \mathbf{y}_j^{(l)}$ is pairwise linearly independent for all i, j , and a.e. choice of $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)}$.

Each node state is updated as the sum

$$\mathbf{x}'_i = \sum_{\mathbf{x}_j \in \mathbb{X}_i} \mathbf{s}_{i,j}$$

of these terms, which is also linearly independent for a.e. choice of $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)}$. Thus, all \mathbf{x}'_i are linearly independent for different $\mathbf{x}_i \in \mathbb{X}$ and $\mathbb{X}_i \subset \mathbb{X}$ which implies injectivity. \square

By this injectivity, the expressivity of l-MCGCs is equivalent to that of the graph isomorphism network (GIN) (Xu et al., 2019), and thus also to the Weisfeiler-Leman test (Weisfeiler & Lehman, 1968). This is the maximal achievable expressivity for models following the message-passing scheme. Due to the universal approximation theorem (Hornik et al., 1989; Hornik, 1991), f can be represented by a multi-layer perceptron (MLP). GIN similarly utilizes an injective function f that is modeled by an MLP, which is applied instead of the linear feature transformation \mathbf{W}_1 with $k = 1$. We will further to this l-MCGC as the multi-channel graph isomorphism network (MC-GIN). This is a critical advantage of MC-GINs: As convolutions are linear operators, it is desired to approximate multi-channel graph convolutions similarly by a linear transformation on \mathbf{X} . While f may introduce a complex transformation, it constructs a linear transformation that is applied to \mathbf{X} . Allowing for learnable edge weights also lifts the constraint of requiring the sum as aggregation.

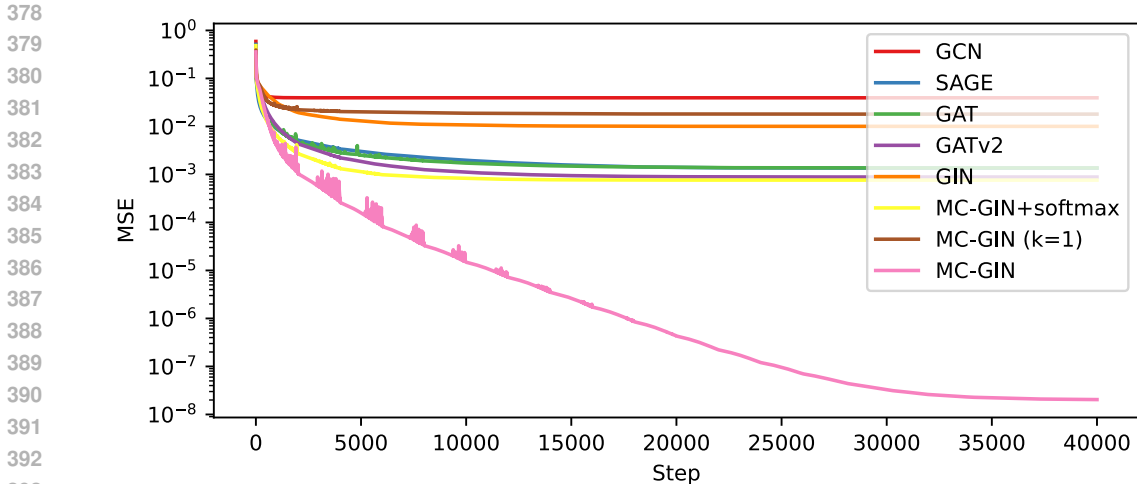


Figure 3: Mean squared error (MSE) during optimization of our function approximation task by applying a single-layer of different methods.

4 EXPERIMENTS

Our aim with this work is to introduce MCGCs in general and to inspire future research to find efficient and effective approximations that will lead to state-of-the-art results. We briefly want to empirically confirm the potential of MCGCs for learning on graphs using one synthetic and one benchmark task. Additional details on models, datasets, and hyperparameters can be found in Appendix A.3. Our implementation is based on PyTorch Geometric (Fey & Lenssen, 2019).

4.1 FUNCTION APPROXIMATION

We first want to evaluate whether the MCGC and the MC-GIN can obtain more informative embeddings by approximating the mapping $\mathbf{Y} = \phi(\mathbf{X})$ where $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times d}$ are multi-channel signals of random but fixed values $X_{i,j} \sim \mathcal{N}(0, 1)$ and $Y_{i,j} \sim \mathcal{N}(0, 1)$. For ϕ , we evaluate a single layer of several graph convolutional operators. Namely the GCN (Kipf & Welling, 2017), GAT (Velickovic et al., 2018), GATv2 (Brody et al., 2022), GIN (Xu et al., 2019), the SAGE convolution (Hamilton et al., 2017), our proposed MC-GIN. We set the number of heads for GAT, GATv2, and MC-GIN to $k = 2$. For functions f in GIN and MC-GIN, we utilize a three-layer MLP with ReLU activations after the first two layers. As an ablation, we further evaluate two versions that do not satisfy our required property for multi-channel MPNNs. We apply a softmax activation to incoming edge weights obtained by f for each node and head with $k = 2$, equivalently to attention (MC-GIN+softmax), and MC-GIN with $k = 1$. We also evaluate the MC-GIN with $k = 3$ and MC-GIN with a one-layer MLP for f and MC-GIN with a two-layer MLP, which satisfy our multi-channel properties. We additionally present results for the MCGC, which, in theory, can represent any such mapping exactly. As the benefits of MCGCs do not depend on specific graph properties, we sample a random undirected Erdős–Rényi graph (Erdős & Rényi, 1959) with n nodes and an edge probability of 5%. We minimize the mean-squared error (MSE) between $\phi(\mathbf{X})$ and \mathbf{Y} using the Adam optimizer. The base learning rate is tuned in $\{0.01, 0.003, 0.001, 0.0003\}$ for each method and is halved every 2000 steps. We set $n = 64$ and $d = 32$.

The averaged loss progression over five runs during optimization is visualized in Figure 3, and the average minimum MSE scores after 40 000 steps are presented in Table 1. MC-GIN achieves an improved approximation error by at least four orders of magnitude compared to all other methods. Notably, all methods that use a single head (GCN, GIN, and MC-GIN ($k=1$)) achieve the worst approximation. Similarly, the error for methods that apply the softmax activation (GAT, GATv2, and MC-GIN+softmax) is significantly higher, as each head amplifies the same signal for all channels. With our ablations, we also find the choice for f not to influence the results strongly (MC-GIN 1 layer MLP, MC-GIN 2 layer MLP) and an additional head to further improve the results (MC-GIN

Method	MSE
GCN	$39 \cdot 10^{-3} \pm 10 \cdot 10^{-4}$
GAT	$14 \cdot 10^{-4} \pm 12 \cdot 10^{-4}$
GATv2	<u>$88 \cdot 10^{-5} \pm 15 \cdot 10^{-4}$</u>
SAGE	$14 \cdot 10^{-4} \pm 46 \cdot 10^{-5}$
GIN	$10 \cdot 10^{-3} \pm 35 \cdot 10^{-4}$
MC-GIN	$20 \cdot 10^{-9} \pm 41 \cdot 10^{-9}$
MC-GIN+softmax	$76 \cdot 10^{-5} \pm 15 \cdot 10^{-4}$
MC-GIN (k=1)	$18 \cdot 10^{-3} \pm 27 \cdot 10^{-4}$
MC-GIN (k=3)	$20 \cdot 10^{-15} \pm 17 \cdot 10^{-15}$
MC-GIN (1 layer MLP)	$18 \cdot 10^{-11} \pm 23 \cdot 10^{-11}$
MC-GIN (2 layer MLP)	$15 \cdot 10^{-9} \pm 17 \cdot 10^{-9}$
MCGC	$37 \cdot 10^{-16} \pm 44 \cdot 10^{-17}$

Table 1: Average and standard deviation of the minimal mean-squared error (MSE). Best MSE in **bold**, second-best underlined.

Method	MAE		Time (s)
	Train	Test	
GCN	6.2 ± 0.2	16.0 ± 0.2	32
GATv2	5.7 ± 0.2	13.6 ± 0.6	63
SAGE	<u>3.9 ± 0.1</u>	<u>12.3 ± 0.2</u>	<u>36</u>
GIN	5.8 ± 0.1	<u>12.3 ± 0.4</u>	42
MC-GIN	3.3 ± 0.2	10.5 ± 0.2	52

Table 2: Results on ZINC. Mean absolute error (MAE) scores are multiplied by 100 for clarity. Best scores in **bold**, second-best underlined.

(k=3). While these improved capabilities come with the risk of overfitting noise in the data, MCGCs are beneficial for complex tasks in which MPNNs struggle to achieve satisfying performance.

4.2 ZINC

We now consider the ZINC dataset (Sterling & Irwin, 2015) to show that approximations of the MCGC can generally also improve benchmark results. It consists of around 250 000 molecular graphs, with the task being to predict the constrained solubility of each molecule. We integrate several base message-passing layers (GCN (Kipf & Welling, 2017), Gatv2 (Brody et al., 2022), SAGE (Hamilton et al., 2017), GIN (Xu et al., 2019), and MC-GIN with $k = 2$) into the implementation of the Long Range Graph Benchmark (Dwivedi et al., 2022) and the updated training scheme of Tönshoff et al. (2024). The number of layers and the learning rate are tuned for each method. As proposed by Dwivedi et al. (2022), each model utilizes at most 500 000 parameters.

Average train and test errors with independent optimal hyperparameters are presented in Table 2. MC-GIN improves both train and test loss by at least 15% compared to the four other methods. The execution time of our implementation is increased by around 24% compared to the GIN. While countless combinations with other techniques, datasets, and tasks can be evaluated, we want to motivate general research on approximating MCGC that will eventually lead to state-of-the-art results.

5 CONCLUSION

In this work, we introduced multi-channel convolutions to graphs by obtaining their form based on the convolution theorem and the graph Fourier transform. While the currently used single-channel graph convolutions, their polynomial approximations, and most MPNNs are not defined for nodes with multiple features, MCGCs are specifically defined for multi-channel signals. This allows MCGCs to obtain more informative embeddings, as different signals can be amplified for each feature channel. Corresponding multi-channel MPNNs need to utilize multiple aggregations and transformations, each amplifying a different signal. We introduce the multi-channel graph isomorphism network (MC-GIN), which can obtain a linear mapping with equivalent expressive power to the graph isomorphism network (Xu et al., 2019). Our experiments confirm the strongly improved abilities to fit complex functions. While our experimental evaluation is limited, the benefits of MCGCs compared to SCGCs are very clear. Having access to the mathematically correct convolution for multi-channel signals allows for the development of various approximations. As more complex interactions between channels may not be required for all tasks, MCGCs increase the risk of overfitting. Approximations of the MCGC will be particularly important as more complex graph-related tasks emerge and more powerful convolutions are needed.

REFERENCES

- 486
487
488 Søren Asmussen. *Applied probability and queues, Second Edition*, volume 51 of *Applications of*
489 *mathematics*. Springer, 2003.
- 490 Michael S. Brandstein and Darren B. Ward (eds.). *Microphone Arrays - Signal Processing*
491 *Techniques and Applications*. Digital Signal Processing. Springer, 2001. doi: 10.1007/
492 978-3-662-04619-7.
- 493 Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? In *The Tenth*
494 *International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29,*
495 *2022*. OpenReview.net, 2022.
- 496 Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally
497 connected networks on graphs. In Yoshua Bengio and Yann LeCun (eds.), *2nd International*
498 *Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,*
499 *Conference Track Proceedings*, 2014.
- 500 John P Burg. Three-dimensional filtering with an array of seismometers. *Geophysics*, 29(5):693–
501 713, 1964.
- 502 Pafnutii Lvovich Chebyshev. *Théorie des mécanismes connus sous le nom de parallélogrammes*.
503 Imprimerie de l’Académie impériale des sciences, 1853.
- 504 Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank
505 graph neural network. In *9th International Conference on Learning Representations, ICLR 2021,*
506 *Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- 507 Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on
508 graphs with fast localized spectral filtering. In Daniel D. Lee, Masashi Sugiyama, Ulrike von
509 Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing*
510 *Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-*
511 *10, 2016, Barcelona, Spain*, pp. 3837–3845, 2016.
- 512 Vijay Prakash Dwivedi, Ladislav Rampásek, Michael Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu,
513 and Dominique Beaini. Long range graph benchmark. In Sanmi Koyejo, S. Mohamed, A. Agar-
514 wal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing*
515 *Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022,*
516 *New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- 517 P Erdős and A Rényi. On random graphs i. *Publicationes Mathematicae Debrecen*, 6:290–297,
518 1959.
- 519 Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric.
520 *CoRR*, abs/1903.02428, 2019. URL <http://arxiv.org/abs/1903.02428>.
- 521 Francesco Di Giovanni, James Rowbottom, Benjamin Paul Chamberlain, Thomas Markovich, and
522 Michael M. Bronstein. Understanding convolution on graphs via energies. *Trans. Mach. Learn.*
523 *Res.*, 2023, 2023.
- 524 William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
525 graphs. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus,
526 S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing*
527 *Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9,*
528 *2017, Long Beach, CA, USA*, pp. 1024–1034, 2017.
- 529 David K Hammond, Pierre Vandergheynst, and Rémi Gribonval. Wavelets on graphs via spectral
530 graph theory. *Applied and Computational Harmonic Analysis*, 30(2):129–150, 2011.
- 531 Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary
532 graph spectral filters via bernstein approximation. In Marc’Aurelio Ranzato, Alina Beygelzimer,
533 Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural In-*
534 *formation Processing Systems 34: Annual Conference on Neural Information Processing Systems*
535 *2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 14239–14251, 2021.

- 540 Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4
541 (2):251–257, 1991. doi: 10.1016/0893-6080(91)90009-T.
- 542 Kurt Hornik, Maxwell B. Stinchcombe, and Halbert White. Multilayer feedforward networks are
543 universal approximators. *Neural Networks*, 2(5):359–366, 1989. doi: 10.1016/0893-6080(89)
544 90020-8.
- 546 Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. OGB-
547 LSC: A large-scale challenge for machine learning on graphs. In Joaquin Vanschoren and Sai-Kit
548 Yeung (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and*
549 *Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021.
- 550 Yinan Huang, William Lu, Joshua Robinson, Yu Yang, Muhan Zhang, Stefanie Jegelka, and Pan
551 Li. On the stability of expressive positional encodings for graphs. In *The Twelfth International*
552 *Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenRe-
553 view.net, 2024.
- 555 Yujiro Inouye and Takehito Sato. Iterative algorithms based on multistage criteria for multichannel
556 blind deconvolution. *IEEE Trans. Signal Process.*, 47(6):1759–1764, 1999. doi: 10.1109/78.
557 765163.
- 558 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net-
559 works. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France,*
560 *April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- 561 Christian Koke and Daniel Cremers. Holonets: Spectral convolutions do extend to directed graphs.
562 In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Aus-*
563 *tria, May 7-11, 2024*. OpenReview.net, 2024.
- 565 Tamara G. Kolda and Brett W. Bader. Tensor decompositions and applications. *SIAM Rev.*, 51(3):
566 455–500, 2009. doi: 10.1137/07070111X.
- 567 Devin Kreuzer, Dominique Beaini, William L. Hamilton, Vincent Létourneau, and Prudencio
568 Tossou. Rethinking graph transformers with spectral attention. In Marc’Aurelio Ranzato, Alina
569 Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in*
570 *Neural Information Processing Systems 34: Annual Conference on Neural Information Process-*
571 *ing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 21618–21629, 2021.
- 572 Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne E.
573 Hubbard, and Lawrence D. Jackel. Handwritten digit recognition with a back-propagation net-
574 work. In David S. Touretzky (ed.), *Advances in Neural Information Processing Systems 2, [NIPS*
575 *Conference, Denver, Colorado, USA, November 27-30, 1989]*, pp. 396–404. Morgan Kaufmann,
576 1989.
- 577 Ron Levie, Federico Monti, Xavier Bresson, and Michael M. Bronstein. Cayleynets: Graph convo-
578 lutional neural networks with complex rational spectral filters. *IEEE Trans. Signal Process.*, 67
579 (1):97–109, 2019. doi: 10.1109/TSP.2018.2879624.
- 581 Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. Gated graph sequence neural
582 networks. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning*
583 *Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceed-*
584 *ings*, 2016.
- 585 Federico Monti, Fabrizio Frasca, Davide Eynard, Damon Mannion, and Michael M. Bronstein. Fake
586 news detection on social media using geometric deep learning. *CoRR*, abs/1902.06673, 2019.
- 587 Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gau-
588 rav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural
589 networks. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019,*
590 *The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The*
591 *Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Hon-*
592 *olulu, Hawaii, USA, January 27 - February 1, 2019*, pp. 4602–4609. AAAI Press, 2019. doi:
593 10.1609/AAAI.V33I01.33014602.

- 594 Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node
595 classification. In *8th International Conference on Learning Representations, ICLR 2020, Addis*
596 *Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- 597
- 598 Richard O’Neil. Convolution operators and $L(p, q)$ spaces. *Duke Mathematical Journal*, 30(1):129
599 – 142, 1963. doi: 10.1215/S0012-7094-63-03015-1.
- 600
- 601 Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Do-
602 minique Beaini. Recipe for a general, powerful, scalable graph transformer. In Sanmi Koyejo,
603 S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural In-*
604 *formation Processing Systems 35: Annual Conference on Neural Information Processing Systems*
605 *2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- 606 Andreas Roth. Simplifying the theory on over-smoothing. *CoRR*, abs/2407.11876, 2024. doi:
607 10.48550/ARXIV.2407.11876.
- 608
- 609 Andreas Roth and Thomas Liebig. Rank collapse causes over-smoothing and over-correlation in
610 graph neural networks. In Soledad Villar and Benjamin Chamberlain (eds.), *Learning on Graphs*
611 *Conference, 27-30 November 2023, Virtual Event*, volume 231 of *Proceedings of Machine Learn-*
612 *ing Research*, pp. 35. PMLR, 2023.
- 613 Andreas Roth, Franka Bause, Nils M Kriege, and Thomas Liebig. Preventing representational rank
614 collapse in mpnns by splitting the computational graph. *CoRR*, abs/2409.11504, 2024.
- 615
- 616 T. Konstantin Rusch, Benjamin Paul Chamberlain, Michael W. Mahoney, Michael M. Bronstein,
617 and Siddhartha Mishra. Gradient gating for deep multi-rate learning on graphs. In *The Eleventh*
618 *International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5,*
619 *2023*. OpenReview.net, 2023.
- 620
- 621 Tara N. Sainath, Ron J. Weiss, Kevin W. Wilson, Bo Li, Arun Narayanan, Ehsan Variani, Michiel
622 Bacchiani, Izhak Shafran, Andrew W. Senior, Kean K. Chin, Ananya Misra, and Chanwoo Kim.
623 Multichannel signal processing with deep neural networks for automatic speech recognition.
624 *IEEE ACM Trans. Audio Speech Lang. Process.*, 25(5):965–979, 2017. doi: 10.1109/TASLP.
625 2017.2672401.
- 626 Michael L Seltzer, Bhiksha Raj, and Richard M Stern. Likelihood-maximizing beamforming for
627 robust hands-free speech recognition. *IEEE Transactions on speech and audio processing*, 12(5):
628 489–498, 2004.
- 629
- 630 Teague Sterling and John J. Irwin. ZINC 15 - ligand discovery for everyone. *J. Chem. Inf. Model.*,
631 55(11):2324–2337, 2015. doi: 10.1021/ACS.JCIM.5B00559.
- 632 Jan Tönshoff, Martin Ritzert, Eran Rosenbluth, and Martin Grohe. Where did the gap go? reassess-
633 ing the long-range graph benchmark. *Trans. Mach. Learn. Res.*, 2024, 2024.
- 634
- 635 Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
636 Bengio. Graph attention networks. In *6th International Conference on Learning Representations,*
637 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.*
638 OpenReview.net, 2018.
- 639
- 640 Mark Weber, Giacomo Domeniconi, Jie Chen, Daniel Karl I. Weidele, Claudio Bellei, Tom Robin-
641 son, and Charles E. Leiserson. Anti-money laundering in bitcoin: Experimenting with graph
642 convolutional networks for financial forensics. *CoRR*, abs/1908.02591, 2019.
- 643 Boris J Weisfeiler and Andrei A Lehman. A reduction of a graph to a canonical form and an algebra
644 arising during this reduction, 1968.
- 645
- 646 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural net-
647 works? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans,*
LA, USA, May 6-9, 2019. OpenReview.net, 2019.

648 Yujun Yan, Milad Hashemi, Kevin Swersky, Yaoqing Yang, and Danai Koutra. Two sides of the
649 same coin: Heterophily and oversmoothing in graph convolutional neural networks. In Xingquan
650 Zhu, Sanjay Ranka, My T. Thai, Takashi Washio, and Xindong Wu (eds.), *IEEE International
651 Conference on Data Mining, ICDM 2022, Orlando, FL, USA, November 28 - Dec. 1, 2022*, pp.
652 1287–1292. IEEE, 2022. doi: 10.1109/ICDM54844.2022.00169.

653 Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. Dive into deep learning. *CoRR*,
654 abs/2106.11342, 2021.

655
656 Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. In *8th International
657 Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
658 OpenReview.net, 2020.

660 A EXPERIMENTAL DETAILS

661
662 Experiments on ZINC were executed on an Nvidia H100 GPU with 96 GB, and experiments on
663 function approximation were executed on an Intel Xeon 8468 Sapphire with 48 cores.

665 A.1 METHODS

666
667 For all methods, we use their standard implementation from PyTorch Geometric (Fey & Lenssen,
668 2019). We mostly use their default parameters but do not add self-loops to any method. For GAT
669 and GATv2, we use two heads and compute the element-wise mean over the outputs of both heads.
670 For GIN, the MLP consists of three linear layers with a bias, followed by a ReLU activation after
671 the first two layers. For MC-GIN, the two adjacent node states of each edge are concatenated, and
672 the same three-layer MLP is applied. Only the number of input channels doubles and the number of
673 output channels is set to the number of heads. The number of heads is set to $k = 2$ unless explicitly
674 stated otherwise.

675 A.2 ZINC

676
677 The ZINC dataset (Sterling & Irwin, 2015) consists of 249 456 graphs, each representing a molecule.
678 Nodes in the graphs represent to heavy atoms, and edges correspond to their connectivity. Each
679 molecule consists of 23.2 nodes and 49.8 edges on average. The task is to predict the constrained
680 solubility, which is a graph regression task. The type of atom is the only feature given. A batch size
681 of 32 is used. The mean absolute error (MAE) is used for optimization and reporting on the test
682 set. Best test scores based on the minimum MAE are considered. Hyperparameters are optimized
683 for each model using a grid search with values for the learning rate in $\{0.001, 0.0003, 0.0001\}$ and
684 the number of layers in $\{1, 2, 4, 8\}$. Each combination is repeated for three random seeds. The
685 minimum average validation score is used to select the hyperparameters for reporting on the test set.
686 The minimum average train scores are directly selected for reporting. ZINC is available under the
687 license DbCL.

688 A.3 HYPERPARAMETERS

689
690 For ZINC, a grid search on the learning rate in $\{0.001, 0.0003, 0.0001\}$ and the number of layers
691 in $\{1, 2, 4, 8\}$ is performed for each method. For the function approximation, a grid search on the
692 learning rate is performed with values $\{0.03, 0.01, 0.003, 0.001\}$. Selected hyperparameters are
693 presented in Table 3 and Table 4.
694
695
696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Method	Loss
GCN	0.01
GAT	0.03
GATv2	0.01
SAGE	0.03
GIN	0.003
MC-GIN+softmax	0.01
MC-GIN (k=1)	0.01
MC-GIN	0.01
MCGC	0.001

Table 3: Best hyperparameters (learning rate (LR)) for Table 1.

Method	Train		Test	
	LR	Layers	LR	Layers
GCN	0.0001	8	0.0001	8
GATv2	0.0003	8	0.0003	8
SAGE	0.0003	8	0.0003	8
GIN	0.0003	8	0.0003	8
MC-GIN	0.0001	8	0.0001	8

Table 4: Best hyperparameters for the results in Table 2. Learning rate (LR) and number of layers (Layers).