

# Graphs, Convolutions, and Neural Networks

## From Graph Filters to Graph Neural Networks

Fernando Gama, Elvin Isufi, Geert Leus, and Alejandro Ribeiro

### Abstract

Network data can be conveniently modeled as a graph signal, where data values are assigned to nodes of a graph that describes the underlying network topology. Successful learning from network data is built upon methods that effectively exploit this graph structure. In this work, we leverage graph signal processing to characterize the representation space of graph neural networks (GNNs). We discuss the role of graph convolutional filters in GNNs and show that any architecture built with such filters has the fundamental properties of permutation equivariance and stability to changes in the topology. These two properties offer insight about the workings of GNNs and help explain their scalability and transferability properties which, coupled with their local and distributed nature, make GNNs powerful tools for learning in physical networks. We also introduce GNN extensions using edge-varying and autoregressive moving average graph filters and discuss their properties. Finally, we study the use of GNNs in recommender systems and learning decentralized controllers for robot swarms.

### Index Terms

Graph signal processing, graph filters, graph convolutions, graph neural networks, stability

## I. INTRODUCTION

Data generated by networks are increasingly common in power grids, robotics, biological, social and economic networks, and recommender systems among others. The irregular and complex

Work in this paper is supported by NSF CCF 1717120, ARO W911NF1710438, ARL DCIST CRA W911NF-17-2-0181, ISTAR-WAS and Intel DevCloud. F. Gama, and A. Ribeiro are with the Dept. of Electrical and Systems Eng., Univ. of Pennsylvania, USA. E. Isufi is with the Multimedia Computing Group and G. Leus is with the Circuits and Systems Group, Delft Univ. of Technology, The Netherlands. E-mails: {fgama, aribeiro}@seas.upenn.edu, {e.isufi-1, g.j.t.leus}@tudelft.nl.

nature of these data poses unique challenges so that successful learning is possible only by incorporating the structure into the inner-working mechanisms of the model [1].

Convolutional neural networks (CNNs) have epitomized the success of leveraging the data structure in temporal series and images transforming the landscape of machine learning in the last decade [2]. CNNs exploit temporal or spatial convolutions to learn an effective nonlinear mapping, scale to large settings, and avoid overfitting [2, Chapter 10]. CNNs offer also some degree of mathematical tractability, allowing to derive theoretical performance bounds under domain perturbations [3]. However, convolutions can only be applied to data in regular domains, hence making CNNs ineffective models when learning from irregular network data.

Graphs are used as a mathematical description of network topologies, while the data can be seen as a signal on top of this graph. In recommender systems, for instance, users can be modeled as nodes, their similarities as edges, and the ratings given to items as graph signals. Processing such data by accounting also for the underlying network structure has been the goal of the field of *graph signal processing* (GSP) [1]. GSP has extended the concepts of Fourier transform, graph convolutions, and graph filtering to process signals while accounting for the underlying topology.

*Graph convolutional* neural networks (GCNNs) build upon graph convolutions to efficiently incorporate the graph structure into the learning process [4]. GCNNs consist of a concatenation of layers, in which each layer applies a *graph convolution* followed by a pointwise nonlinearity [5]–[11]. GCNNs exhibit the key properties of permutation equivariance and stability to perturbations [12], [13]. The former means GCNNs exploit topological symmetries in the underlying graph, while the latter implies the output is robust to small changes in the graph structure. These results allow GCNNs to scale to large graphs and transfer to different (but similar) scenarios.

Graph convolutions can be exactly modeled by finite impulse response (FIR) graph filters [1]. FIR graph filters often require large orders to yield highly discriminatory models, demanding more parameters and an increased computational cost. These limitations are well-understood in the field of GSP and alternative graph filters such as the autoregressive moving average (ARMA) and edge varying graph filters have been proposed to address this [14], [15]. ARMA graph filters maintain the convolutional structure but can achieve a similar response with fewer parameters. Contrarily, the edge varying graph filters are inspired by their time varying counterparts and adapt their structure to the specific graph location. The enhanced flexibility of edge varying graph filters

requires more parameters but their use lays the foundation of a unified framework for all *graph neural networks* (GNNs) [16], generalizing GCNNs by using non-convolutional graph filters.

In this work, we characterize the representation space of GNNs, obtaining properties and insights that hold irrespective of the specific implementation or set of parameters obtained from training. We highlight the role of graph filters in such a characterization and exploit GSP concepts to derive the permutation equivariance and stability properties that hold for all GCNNs. Section II formally introduces graph convolutions. Section III presents the GCNN and discusses permutation equivariance (Sec. III-A) and stability to graph perturbations (Sec. III-B). Section IV generalizes GCNNs by employing alternative graph filters. Section V provides two applications namely rating prediction in recommender systems and learning decentralized controllers for flocking a robot swarm. Section VI contains the paper conclusions and future research directions.

## II. GRAPHS AND CONVOLUTIONS

We capture the irregular structure of the data by means of an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$  with node set  $\mathcal{V} = \{1, \dots, N\}$ , edge set  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  and weight function  $\mathcal{W} : \mathcal{E} \rightarrow \mathbb{R}_+$ . The neighborhood of node  $i \in \mathcal{V}$  is the set of nodes that share an edge with node  $i$  and it is denoted as  $\mathcal{N}_i = \{j \in \mathcal{V} : (j, i) \in \mathcal{E}\}$ . An  $N \times N$  real symmetric matrix  $\mathbf{S}$ , known as the *graph shift operator*, is associated to the graph and satisfies  $[\mathbf{S}]_{ij} = s_{ij} = 0$  if  $(j, i) \notin \mathcal{E}$  for  $j \neq i$ , i.e., the shift operator has a zero whenever two nodes are disconnected. Common shift operators include the adjacency, Laplacian, and Markov matrices as well as their normalized counterparts [1]. The data on top of this graph forms a graph signal  $\mathbf{x} \in \mathbb{R}^N$ , where the  $i$ th entry  $[\mathbf{x}]_i = x_i$  is the datum of node  $i$ . Entries  $x_i$  and  $x_j$  are pairwise related to each other if there exists an edge  $(i, j) \in \mathcal{E}$ . The graph signal  $\mathbf{x}$  can be *shifted* over the nodes by using  $\mathbf{S}$  so that the  $i$ th entry of  $\mathbf{S}\mathbf{x}$  is

$$[\mathbf{S}\mathbf{x}]_i = \sum_{j=1}^N [\mathbf{S}]_{ij} [\mathbf{x}]_j = \sum_{j \in \mathcal{N}_i} s_{ij} x_j. \quad (1)$$

where the last equality holds due to the sparsity of  $\mathbf{S}$  (locality). The output  $\mathbf{S}\mathbf{x}$  is another graph signal where the value at each node is the linear combination of the values of  $\mathbf{x}$  at the neighbors.

Equipped with the notion of signal shift, we define the *graph convolution* as a linear shift-

and-sum operation. Given a set of parameters  $\mathbf{h} = [h_0, \dots, h_K]^\top$ , the graph convolution is

$$\mathbf{H}(\mathbf{S})\mathbf{x} = \sum_{k=0}^K h_k \mathbf{S}^k \mathbf{x}. \quad (2)$$

Operation (2) linearly combines the information contained in different neighborhoods. The  $k$ -shifted signal  $\mathbf{S}^k \mathbf{x}$  contains a summary of the information located in the  $k$ -hop neighborhood and  $h_k$  weighs this summary. This is a *local* operation since  $\mathbf{S}^k \mathbf{x} = \mathbf{S}(\mathbf{S}^{k-1} \mathbf{x})$  entails  $k$  information exchanges with one-hop neighbors [cf. (1)]. The graph convolution (2) *filters* a graph signal  $\mathbf{x}$  with a *FIR graph filter*  $\mathbf{H}(\mathbf{S})$ ; thus, we refer to the weights  $h_k$  as the *filter taps* or *filter weights*.

We can gain additional insight about graph convolutions by analyzing (2) in the *graph frequency domain* [1]. Consider the eigendecomposition of the shift operator  $\mathbf{S} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$  with orthogonal eigenvector matrix  $\mathbf{V} \in \mathbb{R}^{N \times N}$  and diagonal eigenvalue matrix  $\mathbf{\Lambda} \in \mathbb{R}^{N \times N}$  ordered as  $\lambda_1 \leq \dots \leq \lambda_N$ . The eigenvectors  $\mathbf{v}_i$  conform the *graph frequency basis* of graph  $\mathcal{G}$  and can be interpreted as signals representing the *graph oscillating modes*, while the eigenvalues  $\lambda_i$  can be considered as *graph frequencies*. Any signal  $\mathbf{x}$  can be expressed in terms of these graph oscillating modes

$$\tilde{\mathbf{x}} = \mathbf{V}^\top \mathbf{x}. \quad (3)$$

Operation (3) is known as the *graph Fourier transform* (GFT) of  $\mathbf{x}$ , in which entry  $[\tilde{\mathbf{x}}]_i = \tilde{x}_i$  denotes the *Fourier coefficient* associated to graph frequency  $\lambda_i$  and quantifies the contribution of mode  $\mathbf{v}_i$  to the signal  $\mathbf{x}$  [1]. Computing the GFT of the output signal (2) yields

$$\tilde{\mathbf{y}} = \mathbf{V}^\top \mathbf{y} = \sum_{k=0}^K h_k \mathbf{V}^\top \mathbf{V} \mathbf{\Lambda}^k \mathbf{V}^\top \mathbf{x} = \sum_{k=0}^K h_k \mathbf{\Lambda}^k \tilde{\mathbf{x}} = \mathbf{H}(\mathbf{\Lambda}) \tilde{\mathbf{x}} \quad (4)$$

where  $\mathbf{H}(\mathbf{\Lambda})$  is a diagonal matrix with  $i$ th diagonal element  $h(\lambda_i)$  for  $h: \mathbb{R} \rightarrow \mathbb{R}$  given by

$$h(\lambda) = \sum_{k=0}^K h_k \lambda^k. \quad (5)$$

The function in (5) is the *frequency response* of the graph filter  $\mathbf{H}(\mathbf{S})$  and it is determined solely by the filter taps  $\mathbf{h}$ . The effect that a filter has on a signal depends on the specific graph through the instantiation of  $h(\lambda)$  on the eigenvalues  $\lambda_i$  of  $\mathbf{S}$ . It affects the  $i$ th frequency content of  $\tilde{\mathbf{y}}$  as

$$\tilde{y}_i = h(\lambda_i) \tilde{x}_i. \quad (6)$$

That is, the graph convolution (2) modifies the  $i$ th frequency content  $\tilde{x}_i$  of the input signal  $\mathbf{x}$  according to the filter value  $h(\lambda_i)$  at frequency  $\lambda_i$ . Notice the graph convolution is a pointwise operator in the graph frequency domain, in analogy to the convolution in time and images.

### III. GRAPH CONVOLUTIONAL NEURAL NETWORKS

Learning from graph data requires identifying a *representation map*  $\Phi(\cdot)$  between the data  $\mathbf{x}$  and the target representation  $\mathbf{y}$  that leverages the graph structure,  $\mathbf{y} = \Phi(\mathbf{x}; \mathbf{S})$ . The image of  $\Phi$  is known as the *representation space* and determines the space of all possible representations  $\mathbf{y}$  for a given  $\mathbf{S}$  and any input  $\mathbf{x}$ . One example of a representation map is the graph convolution  $\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}) = \mathbf{H}(\mathbf{S})\mathbf{x}$  in (2), where set  $\mathcal{H} = \{\mathbf{h}\}$  contains the filter coefficients that characterize its representation space [17]. To *learn* this map, we consider a cost function  $J(\cdot)$  and a training set  $\mathcal{T} = \{\mathbf{x}_1, \dots, \mathbf{x}_{|\mathcal{T}|}\}$  with  $|\mathcal{T}|$  samples. The *learned map* is then  $\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}^*)$  with

$$\mathcal{H}^* = \underset{\mathcal{H}}{\operatorname{argmin}} \frac{1}{|\mathcal{T}|} \sum_{\mathbf{x} \in \mathcal{T}} J(\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})). \quad (7)$$

Typical cost functions include the mean squared error (MSE) or the L1 loss for regression and cross-entropy loss for classification problems [2]. Problem (7) consists of finding the  $K+1$  filter taps  $\mathcal{H}^* = \{\mathbf{h}^*\}$  that best fit the training data w.r.t. cost  $J(\cdot)$ , with  $K$  being a design choice (a hyperparameter). However, graph convolutions limit the representation power to linear mappings. We can increase the class of mappings that leverage the graph by nesting convolutions into a nonlinearity. The latter leads to the concept of *graph perceptron*, which is formalized next.

**Definition 1** (Graph perceptron). A graph perceptron is a mapping that applies an entrywise nonlinearity  $\sigma(\cdot)$  to the output of a graph convolution  $\mathbf{H}(\mathbf{S})\mathbf{x}$ , i.e.,

$$\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}) = \sigma(\mathbf{H}(\mathbf{S})\mathbf{x}), \quad (8)$$

where set  $\mathcal{H} = \{\mathbf{h}\}$  contains the filter coefficients.

The graph perceptron generates another graph signal obtained as a graph convolution followed by a nonlinearity (e.g., a rectified linear unit, or ReLU for short,  $\sigma(z) = \max\{z, 0\}$ ). As such, the graph perceptron captures nonlinear relationships between the data  $\mathbf{x}$  and the target representation  $\mathbf{y}$ . By building then a cascade of  $L$  graph perceptrons, we get a *multi-layer graph perceptron*,

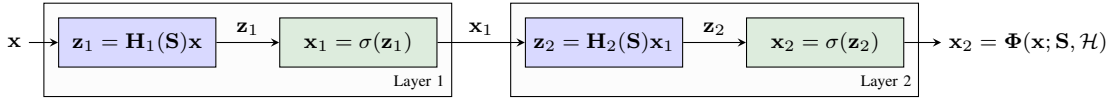


Figure 1. Each blue block represents a linear graph filter and each green block represents a nonlinearity. The concatenation of a convolutional graph filter [cf. (2)] and a nonlinearity forms a graph perceptron (Def. 1) or *layer*. Using a bank of graph convolutional filters [cf. (10)-(11)] and cascading several layers, leads to a GCNN. GCNNs are a subset of GNNs, which follow the same structure but consider arbitrary graph filters, see Section IV.

where at layer  $\ell$  we compute

$$\mathbf{x}_\ell = \sigma(\mathbf{H}_\ell(\mathbf{S})\mathbf{x}_{\ell-1}), \quad \ell = 1, \dots, L. \quad (9)$$

Differently from (8), a multi-layer graph perceptron allows nonlinear signal mixing between nodes. This can be seen in (9) where the input of the perceptron at layer  $\ell$  is  $\mathbf{x}_{\ell-1}$ , which is in turn the output of the perceptron at layer  $\ell - 1$ ,  $\mathbf{x}_{\ell-1} = \sigma(\mathbf{H}_{\ell-1}(\mathbf{S})\mathbf{x}_{\ell-2})$ . The cascade form allows, therefore, graph convolutions of nonlinear signal transformations coming from the precedent layer. Unrolling this recursion to all layers, we have that the input to the first layer is the data  $\mathbf{x}_0 = \mathbf{x}$  and the output of the last layer is the estimate of the target representation  $\mathbf{x}_L = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$ ; here, the set  $\mathcal{H} = \{\mathbf{h}_\ell\}_\ell$  contains the filter taps of the  $L$  graph filters in (9).

The graph perceptron (8) and the multi-layer graph perceptron (9) can be viewed as specific *graph convolutional neural networks* (GCNNs). The former is a GCNN of one layer, while the latter is a GCNN of  $L$  layers. As it is a good practice in neural networks [2], we can substantially increase the representation power of GCNNs by incorporating multiple parallel features per layer. These features are the result of processing multiple input features with a parallel bank of graph filters. Let us consider  $F_{\ell-1}$  input graph signal features  $\mathbf{x}_{\ell-1}^1, \dots, \mathbf{x}_{\ell-1}^{F_{\ell-1}}$  at layer  $\ell$ . Each input feature  $\mathbf{x}_{\ell-1}^g$  for  $g = 1, \dots, F_{\ell-1}$  is processed in parallel by  $F_\ell$  different graph filters of the form (2) to output the  $F_\ell$  convolutional features

$$\mathbf{u}_\ell^{fg} = \mathbf{H}_\ell^{fg}(\mathbf{S})\mathbf{x}_{\ell-1}^g = \sum_{k=0}^K h_{\ell k}^{fg} \mathbf{S}^k \mathbf{x}_{\ell-1}^g, \quad f = 1, \dots, F_\ell. \quad (10)$$

The convolutional features are subsequently summarized along the input index  $g$  to yield the aggregated features (see [16, eq. (13)] for a compact matrix-based notation)

$$\mathbf{u}_\ell^f = \sum_{g=1}^{F_{\ell-1}} \mathbf{H}_\ell^{fg}(\mathbf{S})\mathbf{x}_{\ell-1}^g, \quad f = 1, \dots, F_\ell. \quad (11)$$

The aggregated features are finally passed through a nonlinearity to complete the  $\ell$ th layer output

$$\mathbf{x}_\ell^f = \sigma(\mathbf{u}_\ell^f), \quad f = 1, \dots, F_\ell. \quad (12)$$

A GCNN in its complete form<sup>1</sup> is a concatenation of  $L$  layers, in which each layer computes operations (10)-(11)-(12). Differently from the multi-layer graph perceptron GCNN in (9), the complete form employs a parallel bank of  $F_\ell \times F_{\ell-1}$  graph convolutional filters. This increases the representation power of the mapping and exploits both the stable operation in signal processing, the *convolution*, and the underlying graph structure of the data. The input to the first layer is the data  $\mathbf{x}_0 = \mathbf{x}$  and the target representation is the collection of  $F_L$  features of the last layer  $[\mathbf{x}_L^1, \dots, \mathbf{x}_L^{F_L}] = \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$ , where set  $\mathcal{H} = \{\mathbf{h}_\ell^{fg}\}_{\ell fg}$  collects now the filter taps of all layers. For a given  $\mathbf{S}$  and fixed hyperparameters  $L$ ,  $F_\ell$  and  $K$ , the representation space of a GCNN is characterized by the set of filter coefficients at each layer  $\mathcal{H}$ . This representation space is different from the one obtained by using linear FIR filters representation maps [17].

We can *learn* the filter taps by solving problem (7) with the GCNN map  $\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H})$ . To do so, we use some optimization method based on stochastic gradient descent [18] and, noting the GCNN is a compositional layered architecture, we also use backpropagation to compute the derivatives of the loss function  $J(\cdot)$  with respect to the filter taps  $\mathcal{H}$  [19]. Since the training data comes from a distribution that has a graph structure  $\mathbf{S}$ , we expect the learned map  $\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}^*)$  to *generalize* and perform well for data  $\mathbf{x} \notin \mathcal{T}$  that come from a similar distribution leveraging  $\mathbf{S}$ . The rationale behind this expectation is that the GCNN is a nonlinear processing architecture that exploits the knowledge the graph carries about the data. Another advantage of a GCNN is its local implementation due to the use of graph convolutions [cf. (2)] and pointwise nonlinearities. In fact, all the  $F_\ell \times F_{\ell-1}$  convolutional features in (10) are local over the graph as they simply comprise a parallel bank of graph convolutional filters, each of which is local [cf. (2)]. Further, since the aggregation step in (11) happens across features of the same node and the nonlinearity in (12) is pointwise, these operations are also local and distributable. This built-in characteristic of GCNNs naturally leads to learning solutions that are distributed on the underlying graph.

<sup>1</sup>We omit pooling to emphasize the role of graph filters. Please, refer to [5]–[7] for pooling methods.

*Implementations of GCNNs.* Given a matrix representation  $\mathbf{S}$  and fixed set of hyperparameters (number of layers  $L$ , filter taps  $K$ , features  $F_\ell$  and nonlinearity  $\sigma(\cdot)$ ), the representation space of the GCNN model (10)-(11)-(12) is characterized by the set of parameters  $\mathcal{H}$  that determine the graph filters. There exist in the literature different implementations for the graph convolution operation (10), as well as other parametrizations that further restrict this representation space. We overview these in light of the description (10)-(11)-(12).

**Same representation space.** Spectral GCNNs [5] compute (10) in the spectral domain (4) and consider the (normalized) Laplacian as the shift  $\mathbf{S}$ ; as long as all the eigenvalues of  $\mathbf{S}$  are different, both (4) and (10) are equivalent [1]. ChebNets [6] use a Chebyshev polynomial to compute the graph convolution and consider as  $\mathbf{S}$  a normalized version of the Laplacian that forces all eigenvalues to be in  $[-1, 1]$  which is required for the use of Chebyshev polynomials; Chebyshev polynomials are equivalent to the polynomials in (10). In summary, we see that [5], [6] just differ in their implementation of the graph convolution, but all cover the same representation space as the GCNN model (10)-(11)-(12) for the specific shift  $\mathbf{S}$ .

**Smaller representation space.** GCNs [8] consider (10) with only the one-hop filter tap  $h_{\ell 1}^{fg}$  for each layer and each filter, i.e.  $K = 1$  and  $h_{\ell 0}^{fg} = 0$  for all  $\ell$ ; they adopt a normalized self-looped version of the adjacency as  $\mathbf{S}$ . Simple graph convolutional networks (SGCs) [9] consider (10) with only the  $K$ -hop filter tap, i.e.  $h_{\ell k}^{fg} = 0$  for all  $k < K$ ; they also adopt a normalized self-looped version of the adjacency as  $\mathbf{S}$ . Graph isomorphism networks (GINs) [10] consider an order-one polynomial  $K = 1$  but with  $h_{\ell 0}^{fg} = (1 + \varepsilon_\ell)h_{\ell 1}^{fg}$  for some pre-defined  $\varepsilon_\ell$ ; it adopts the binary adjacency as  $\mathbf{S}$  and suggests the inclusion of layers with  $K = 0$  in between layers with  $K = 1$ . Diffusion CNNs [11] consider a single layer with  $F_1 = NF_0$  and the same  $K$  filter taps for all input features  $h_{1k}^{fg} = h_{1k}^f$ ; it adopts the adjacency matrix as  $\mathbf{S}$ . It follows that the representations space of [8]–[11] is just a subspace of the representation space of the GCNN model in (10)-(11)-(12).

We note that, while the representation space of [5], [6] is the same as in the GCNN model (10)-(11)-(12), their difference in the implementation of the graph convolution impacts how the optimization space is navigated during training, arriving at different solutions. No particular implementation, however, has consistently outperformed the rest across a wide range of problems. In any case, since the representation space is the same, the characterizations, properties and insights established here apply to all of these. Implementations [8]–[11], on the other hand, further regularize the graph convolution, constraining the representation space to be a subspace of that in the GCNN model. These might be useful in problems with smaller datasets, or where further information on the data structure is available.



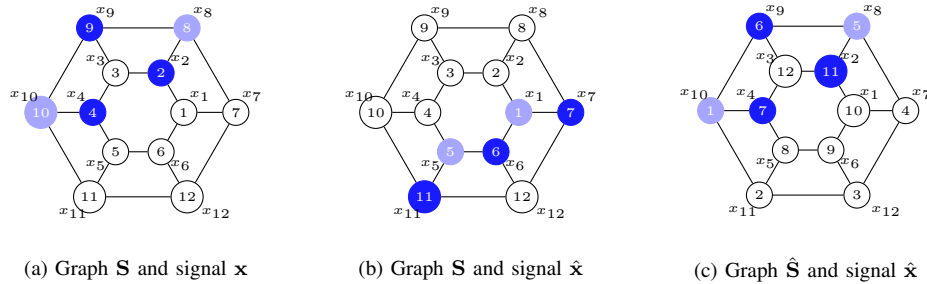


Figure 2. Permutation equivariance of GCNNs. The output of a GCNN is equivariant to graph permutations (Theorem 1). This means independence from labeling and shows GCNNs exploit internal signal symmetries. Signals in (a) and (b) are different on the same graph but they are permutations of each other –interchange inner and outer hexagons and rotate  $180^\circ$  [c.f. (c)]. A GCNN would learn how to classify the signal in (b) from seeing examples of the signal in (a). Integers represent labels, while colors signal values.

### A. Permutation equivariance

A graph shift operator  $\mathbf{S}$  fixes an arbitrary ordering of the nodes in the graph. Since nodes are naturally unordered, we want the GCNN output to be unaffected by it. That is, we want any change in node ordering to be reflected with the corresponding reordering in the GCNN output. It turns out GCNNs are unaffected by node labeling –a property known as permutation equivariance– as stated by the following theorem.

**Theorem 1** (Permutation Equivariance [12], [13]). Consider an  $N \times N$  permutation matrix  $\mathbf{P}$  and the permutations of the shift operator  $\hat{\mathbf{S}} = \mathbf{P}^\top \mathbf{S} \mathbf{P}$  and of the input data  $\hat{\mathbf{x}} = \mathbf{P}^\top \mathbf{x}$ . For a GCNN  $\Phi(\cdot)$ , it holds that

$$\Phi(\hat{\mathbf{x}}; \hat{\mathbf{S}}, \mathcal{H}) = \mathbf{P}^\top \Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}). \quad (13)$$

Theorem 1 states that a node reordering results in a corresponding reordering of the GCNN output, implying GCNNs are independent of node labeling. Theorem 1 implies also that graph convolutions exploit the inherent symmetries present in a graph to improve data processing. If the graph exhibits several nodes with the same topological neighborhood (graph symmetries), then learning how to process data in any of these nodes can be translated to every other node with the same topological neighborhood. This allows GCNNs to learn from fewer samples and generalize easier to signals located at any topologically similar neighborhood, see Fig. 2.

### B. Stability to perturbations

Since real graphs rarely exhibit perfect symmetries, we are interested in more general changes to the underlying graph support than just permutations. For instance, in problems where the

graph  $\mathbf{S}$  is fixed but unknown, we need to use an estimate  $\hat{\mathbf{S}}$  of it but we still want the GCNN to work well as long as the estimate is good (Section V-A). On another set of problems, the graph support may naturally differ from training  $\mathbf{S}$  to testing  $\hat{\mathbf{S}}$ , a scenario known as *transfer learning* (Section V-B). In these cases, we need the GCNN to have a similar performance whether they run on  $\mathbf{S}$  or on  $\hat{\mathbf{S}}$  as long as both graphs are similar. To measure the similarity between graphs  $\mathbf{S}$  and  $\hat{\mathbf{S}}$ , and in light of Theorem 1, we define next the relative distance modulo permutation.

**Definition 2** (Relative distance). Consider the set of all permutation matrices  $\mathcal{P} = \{\mathbf{P} \in \{0, 1\}^{N \times N} : \mathbf{P}^T \mathbf{1} = \mathbf{1}, \mathbf{P} \mathbf{1} = \mathbf{1}\}$ . For two graphs  $\mathbf{S}$  and  $\hat{\mathbf{S}}$  with the same number of nodes, we define the set of *relative error matrices* as

$$\mathcal{R}(\mathbf{S}, \hat{\mathbf{S}}) = \{\mathbf{E} : \mathbf{P}^T \hat{\mathbf{S}} \mathbf{P} = \mathbf{S} + (\mathbf{E} \mathbf{S} + \mathbf{S} \mathbf{E}), \mathbf{P} \in \mathcal{P}\}. \quad (14)$$

The *relative distance* modulo permutations between  $\mathbf{S}$  and  $\hat{\mathbf{S}}$  is then defined as

$$d(\mathbf{S}, \hat{\mathbf{S}}) = \min_{\mathbf{E} \in \mathcal{R}(\mathbf{S}, \hat{\mathbf{S}})} \|\mathbf{E}\| \quad (15)$$

where  $\|\cdot\|$  indicates the operator norm. We denote by  $\mathbf{E}^*$  and  $\mathbf{P}^*$  the relative error matrix and the permutation matrix that minimize (15), respectively.

We readily see that if  $d(\mathbf{S}, \hat{\mathbf{S}}) = 0$ , then  $\hat{\mathbf{S}}$  is a permutation of  $\mathbf{S}$ , and thus the relative distance (15) measures how far  $\mathbf{S}$  and  $\hat{\mathbf{S}}$  are from being permutations of each other. We note that, unlike the absolute perturbation model, the relative distance (Def. 2) accurately reflects changes to both the edge weights and the topology structure [12].

The change in the output of a GCNN due to a change in the underlying support is bounded for GCNNs whose constitutive graph filters are integral Lipschitz.

**Definition 3** (Integral Lipschitz filters). We say a filter  $\mathbf{H}(\mathbf{S})$  is *integral Lipschitz* if its frequency response  $h(\lambda)$  [cf. (5)] is such that  $|h(\lambda)| \leq 1$  and its derivative  $h'(\lambda)$  satisfies  $|\lambda h'(\lambda)| \leq C$  for some finite constant  $C$ .

The derivative condition  $|\lambda h'(\lambda)| \leq C$  implies integral Lipschitz filters have frequency responses that can vary rapidly around  $\lambda = 0$  but are flat for  $\lambda \rightarrow \infty$ , see Fig. 3a. GCNNs that use integral Lipschitz filters are stable under relative perturbations. This means the change in the output due to changes in the underlying graph is bounded by the size of the perturbation [cf. Def. 2].

**Theorem 2** (Stability [12]). Let  $\mathbf{S}$  and  $\hat{\mathbf{S}}$  be two different graphs with the same number of nodes such that their relative distance is  $d(\mathbf{S}, \hat{\mathbf{S}}) \leq \varepsilon$  [cf. Def. 2]. Let  $\Phi(\cdot; \cdot, \mathcal{H})$  be a multi-layer graph perceptron GCNN [cf. (9)] where all filters  $\mathcal{H}$  are integral Lipschitz with constant  $C$  [cf. Def. 3]. Then, it holds that

$$\|\Phi(\mathbf{x}; \mathbf{S}, \mathcal{H}) - \Phi(\mathbf{P}^{\star\top} \mathbf{x}; \mathbf{P}^{\star\top} \hat{\mathbf{S}} \mathbf{P}^{\star}, \mathcal{H})\| \leq 2C (1 + \delta\sqrt{N}) L \varepsilon \|\mathbf{x}\| + \mathcal{O}(\varepsilon^2) \quad (16)$$

where  $\delta = (\|\mathbf{U} - \mathbf{V}\| + 1)^2 - 1$  is the eigenvector misalignment between the eigenbasis  $\mathbf{V}$  of  $\mathbf{S}$  and the eigenbasis  $\mathbf{U}$  of the relative error matrix  $\mathbf{E}^*$ , with  $\mathbf{E}^*$  and  $\mathbf{P}^*$  given in Definition 2.

Theorem 2 proves that a change  $\varepsilon$  in the shift operator causes a change proportional to  $\varepsilon$  in the GCNN output. The proportionality constant has the term  $C$  that depends on the filter design, and the term  $(1 + \delta\sqrt{N})$  that depends on the specific perturbation. But it also has a constant factor  $L$  that depends on the depth of the architecture implying deeper GCNNs are less stable.

#### IV. EXTENSIONS: GENERAL GRAPH FILTERS

Oftentimes, the GCNN would require highly sharp filter responses to discriminate between classes. We can increase the discriminatory power by either increasing the filter order  $K$  or changing the filter type  $\mathbf{H}(\mathbf{S})$  in the graph perceptron (8). Increasing  $K$  is not always feasible as it leads to more filter coefficients, a higher complexity, and numerical issues related to the higher order powers of the shift operator  $\mathbf{S}^k$ . Instead, changing the filter type allows implementing another family of graph neural networks (GNNs) with different properties. We present two alternative filters that provide different insights on how to design more general GNNs: *the autoregressive moving average* (ARMA) graph filter [14] and *the edge varying* graph filter [15].

##### A. ARMA $Net$

An ARMA graph filter operates also pointwise in the spectral domain  $\tilde{y}_i = h(\lambda_i)\tilde{x}_i$  [cf. (6)] but it is characterized by the rational frequency response

$$h(\lambda) = \frac{\sum_{q=0}^Q b_q \lambda^q}{1 + \sum_{p=1}^P a_p \lambda^p}. \quad (17)$$

The frequency response is now controlled by  $P$  denominator coefficients  $\mathbf{a} = [a_1, \dots, a_P]^\top$  and  $Q + 1$  numerator coefficients  $\mathbf{b} = [b_0, \dots, b_Q]^\top$ . The rational frequency responses in (17) span

*Insights on stability.* To offer further insight into Theorem 2, consider the particular case where the perturbation  $\hat{\mathbf{S}}$  is an edge dilation of a graph  $\mathbf{S}$ , i.e.  $\hat{\mathbf{S}} = (1 + \varepsilon)\mathbf{S}$ , where all edges are increased proportionally by a factor  $\varepsilon$ . The relative error matrix is  $\mathbf{E} = (\varepsilon/2)\mathbf{I}$  so that the relative distance is  $d(\mathbf{S}, \hat{\mathbf{S}}) = \|\mathbf{E}\| \leq \varepsilon$ . The graph dilation changes the eigenvalues to  $\hat{\lambda}_i = (1 + \varepsilon)\lambda_i$  while the eigenvectors remain the same. We note that, even if  $\varepsilon$  is small, the change in eigenvalues could be large if  $\lambda_i$  is large, see Fig. 3a.

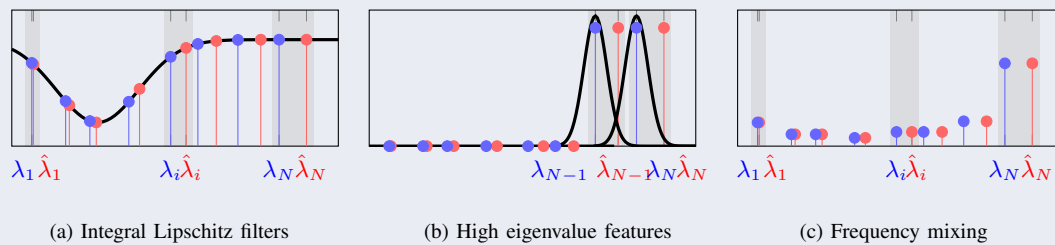


Figure 3. (a) Frequency response for an integral Lipschitz filter (in black), eigenvalues for  $\mathbf{S}$  (in blue) and eigenvalues for  $\hat{\mathbf{S}}$  (in red). Larger eigenvalues exhibit a larger change. (b) Separating energy located at  $\lambda_{N-1}$  from that at  $\lambda_N$  requires filters with sharp transitions that are not integral Lipschitz. Then, a change in eigenvalues renders these filters useless (they are not stable). (c) Applying a ReLU to a signal with all its energy located at  $\lambda_N$  results in a signal with energy spread through the spectrum. Information on low eigenvalues can be discriminated in a stable fashion.

This observation that even small perturbations lead to large changes in the eigenvalues can considerably affect the output of a graph filter causing instability, unless the graph filters are carefully designed. To see this, consider first the output of a graph filter in the frequency domain,  $\tilde{y}_i = h(\lambda_i)\tilde{x}_i$  [cf. (6)]. With the graph dilation, the frequency response gets instantiated at  $\hat{\lambda}_i = (1 + \varepsilon)\lambda_i$  instead of  $\lambda_i$ , so the  $i$ th frequency content is now  $\hat{y}_i = h(\hat{\lambda}_i)\tilde{x}_i$ . The change between the original  $i$ th frequency content of the output  $\tilde{y}_i$  and the perturbed one  $\hat{y}_i$  depends on how much  $h(\lambda_i)$  changes with respect to  $h(\hat{\lambda}_i)$ , and thus can be quite large for large  $\lambda$ . So if we want  $\tilde{y}_i$  to be close to  $\hat{y}_i$  for stability, we need to have frequency responses  $h(\lambda)$  that have a flat response for large eigenvalues, see Fig. 3a. Integral Lipschitz filters do have a flat response at large eigenvalues and thus are stable.

The cost to pay for stability is that integral Lipschitz filters cannot discriminate information located at higher eigenvalues. As seen in Fig. 3b, discriminative filters are narrow filters. Then, if even a small perturbation causes a large change in the instantiated eigenvalue (as is the case for large eigenvalues), the filter output changes to a zero output, and thus is not stable. Thus, linear graph filters exhibit a trade-off between discriminability and stability; a trait shared by regular convolutional filters [3].

GCNNs incorporate pointwise nonlinearities in the graph perceptron. This nonlinear operation has a frequency mixing effect (akin to demodulation) by which the signal energy is spilled throughout the spectrum, see Fig. 3c. Thus, energy from large eigenvalues now appears in smaller eigenvalues. This new low-eigenvalue frequency content can be captured and discriminated by subsequent filters in a stable manner. Therefore, pointwise nonlinearities make GCNNs information processing architectures that are both stable and selective.

an equivalent space to that of graph filters in (2). However, the spectral equivalence does not imply that the two filters have the same properties. ARMA filters implement rational frequency responses rather than polynomial ones as FIR filters do (5). Therefore, we expect them to achieve a sharper response with lower orders of  $P$  and  $Q$  such that  $P + Q < K$ . Replacing the spectral variable  $\lambda$  with the shift operator  $\mathbf{S}$  allows us to write the ARMA output  $\mathbf{y} = \mathbf{H}(\mathbf{S})\mathbf{x}$  as

$$\mathbf{y} = \left( \mathbf{I} + \sum_{p=1}^P a_p \mathbf{S}^p \right)^{-1} \left( \sum_{q=0}^Q b_q \mathbf{S}^q \right) \mathbf{x} := \mathbf{P}(\mathbf{S})^{-1} \mathbf{Q}(\mathbf{S}) \mathbf{x} \quad (18)$$

where  $\mathbf{P}(\mathbf{S}) := \mathbf{I} + \sum_{p=1}^P a_p \mathbf{S}^p$  and  $\mathbf{Q} := \sum_{q=0}^Q b_q \mathbf{S}^q$  are two FIR filters [cf. (2)] that allow writing the ARMA filter as  $\mathbf{H}(\mathbf{S}) = \mathbf{P}(\mathbf{S})^{-1} \mathbf{Q}(\mathbf{S})$ . As it follows from (18), we need to apply the matrix inverse  $\mathbf{P}(\mathbf{S})$  to obtain the ARMA output. This, unless the number of nodes is moderate, is computationally unaffordable; hence, we need an iterative method to approximately apply the inverse. Due to its faster convergence, we choose a parallel structure that consists of first transforming the polynomial ratio in (18) into its partial fraction decomposition form and subsequently using the Jacobi method to approximately apply the inverse. While also other Krylov approaches are possible to solve (18), the parallel Jacobi method offers a better tradeoff between computational complexity, distributed implementation, and convergence.

**Partial fraction decomposition of ARMA filters.** Consider the rational frequency response  $h(\lambda)$  in (17) and let  $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_P]^\top$  be the  $P$  poles,  $\boldsymbol{\beta} = [\beta_1, \dots, \beta_P]^\top$  the corresponding residuals and  $\boldsymbol{\alpha} = [\alpha_0, \dots, \alpha_K]^\top$  the direct terms. Then, we can write (18) in the equivalent form

$$\mathbf{y} = \sum_{p=1}^P \beta_p \left( \mathbf{S} - \gamma_p \mathbf{I} \right)^{-1} \mathbf{x} + \sum_{k=0}^K \alpha_k \mathbf{S}^k \mathbf{x}. \quad (19)$$

The equivalence of (19) and (18) implies that instead of learning  $\mathbf{a}$  and  $\mathbf{b}$  in (18), we can learn  $\boldsymbol{\alpha}$ ,  $\boldsymbol{\beta}$ , and  $\boldsymbol{\gamma}$  in (19). To avoid the matrix inverses in the single pole filters, we can approximate each output  $\mathbf{u}_p$  through the Jacobi method.

**Jacobi method for single pole filters.** We can write the output of the  $p$ th single pole filter  $\mathbf{u}_p$  in the equivalent linear equation form  $(\mathbf{S} - \gamma_p \mathbf{I}) \mathbf{u}_p = \beta_p \mathbf{x}$ . The Jacobi algorithm requires separating  $(\mathbf{S} - \gamma_p \mathbf{I})$  into its diagonal and off-diagonal terms. Defining  $\mathbf{D} = \text{diag}(\mathbf{S})$  as the matrix containing the diagonal of the shift operator, we can write the Jacobi approximation  $\mathbf{u}_{p\tau}$  of the  $p$ th single

pole filter output  $\mathbf{u}_p$  at iteration  $\tau$  by the recursive expression

$$\mathbf{u}_{p\tau} = \left(\mathbf{D} - \gamma_p \mathbf{I}\right)^{-1} \left[ \beta_p \mathbf{x} - \left(\mathbf{S} - \mathbf{D}\right) \mathbf{u}_{p(\tau-1)} \right], \quad \text{with } \mathbf{u}_{p0} = \mathbf{x}. \quad (20)$$

The inverse in (20) is now element-wise on the diagonal matrix  $(\mathbf{D} - \gamma_p \mathbf{I})$ . This recursion can be unrolled to all its terms to write an explicit relationship between  $\mathbf{u}_{pT}$  and  $\mathbf{x}$ . To do that, we define the parameterized shift operator  $\mathbf{R}(\gamma_p) = -(\mathbf{D} - \gamma_p \mathbf{I})^{-1}(\mathbf{S} - \mathbf{D})$  and use it to write the  $T$ th Jacobi recursion as

$$\mathbf{u}_{pT} = \beta_p \sum_{\tau=0}^{T-1} \mathbf{R}^\tau(\gamma_p) \mathbf{x} + \mathbf{R}^T(\gamma_p) \mathbf{x}. \quad (21)$$

For a convergent Jacobi method,  $\mathbf{u}_{pT}$  converges to the single pole output  $\mathbf{u}_p$ . However, in a practical setting we truncate (21) for a finite  $T$ . We can then write the single pole filter output as  $\mathbf{u}_{pT} := \mathbf{H}_T(\mathbf{R}(\gamma_p))\mathbf{x}$ , where we define the following FIR filter of order  $T$

$$\mathbf{H}_T(\mathbf{R}(\gamma_p)) = \beta_p \sum_{\tau=0}^{T-1} \mathbf{R}^\tau(\gamma_p) + \mathbf{R}^T(\gamma_p). \quad (22)$$

with the parametric shift operator  $\mathbf{R}(\gamma)$ . In other words, a single pole filter is approximated by a graph convolutional filter of the form (2) in which the shift operator  $\mathbf{S}$  is substituted by  $\mathbf{R}(\gamma)$ . This parametric convolutional filter uses coefficients  $\beta_p$  for  $\tau = 0, \dots, T-1$  and 1 for  $\tau = T$ .

**Jacobi ARMA filters and ARMANets.** Assuming we use truncated Jacobi iterations of order  $T$  to approximate all single pole filters in (19), we can write the ARMA filter as

$$\mathbf{H}(\mathbf{S}) = \sum_{p=1}^P \mathbf{H}_T(\mathbf{R}(\gamma_p)) + \sum_{k=0}^K \alpha_k \mathbf{S}^k \quad (23)$$

where the  $p$ th approximated single pole filter  $\mathbf{H}_T(\mathbf{R}(\gamma_p))$  is defined in (22) and the parametric shift operator  $\mathbf{R}(\gamma_p)$  in (21). In summary, a Jacobi approximation of the ARMA filter with orders  $(P, T, K)$  is the one defined by (22) and (23). Scalar  $P$  indicates the number of poles,  $T$  the number of Jacobi iterations, and  $K$  the order of the direct term  $\sum_{k=0}^K \alpha_k \mathbf{S}^k$  in (19).

Substituting (23) into (8) yields an ARMA graph perceptron, which is the building block for ARMA GNNs or, for short, ARMANets. ARMANets are themselves convolutional. For a sufficiently large number of Jacobi iterations  $T$ , (23) is equivalent to (18) which performs a pointwise multiplication in the spectral domain with the response (17). The Jacobi filters in (23) are also reminiscent of the convolutional filters in (2). But the similarity is superficial because

in ARMANets we train also the  $2P$  single pole filter coefficients  $\beta_p$  and  $\gamma_p$  alongside the  $K + 1$  coefficients of the direct term  $\sum_{k=0}^K \alpha_k \mathbf{S}^k$ . The equivalence suggests ARMANets may help achieve more discriminatory filters by tuning the single pole filter orders  $P$  and  $T$ . An example of an implementation of ARMANets are CayleyNets [20], see [16].

### B. EdgeNet

While ARMANets enhance the discriminatory power of GCNNs with alternative convolutional filters, the edge varying GNN departs from the convolutional prior to improve GCNNs. The EdgeNet leverages the sparsity and locality of the shift operator  $\mathbf{S}$  and forms a graph perceptron [cf. (9)] by replacing the graph convolutional filter with an edge varying graph filter [15].

**From shared to edge parameters.** In the convolutional filter (2), all nodes share the same scalar  $h_k$  to weigh equally the information from all  $k$ -hop away neighbors  $\mathbf{S}^k \mathbf{x}$ . This is advantageous because it limits the number of trainable parameters, allows permutation equivariance, and favors stability. However, this parameter sharing limits also the discriminatory power to architectures whose filters  $\mathbf{H}(\mathbf{S})$  have the same eigenvectors as  $\mathbf{S}$  [cf. (4)]. We can improve the discriminatory power by considering a linear filter in which node  $i$  uses a scalar  $\Phi_{ij}^{(k)}$  to weigh the information of its neighbor  $j$  at iteration  $k$ . For  $k = 0$ , each node weighs only its own signal to build the zero-shifted signal  $\mathbf{z}^{(0)} = \Phi^{(0)} \mathbf{x}$ , where  $\Phi^{(0)}$  is an  $N \times N$  diagonal matrix of parameters with  $i$ th diagonal entry  $\Phi_{ii}^{(0)}$  being the weight of node  $i$ . Signal  $\mathbf{z}^{(0)}$  is subsequently exchanged with neighboring nodes to build the one-shifted signal  $\mathbf{z}^{(1)} = \Phi^{(1)} \mathbf{z}^{(0)}$ , where the parameter matrix  $\Phi^{(1)}$  shares the support with  $\mathbf{I} + \mathbf{S}$ ; the  $(i, j)$ th entry  $\Phi_{ij}^{(1)}$  is the weight node  $i$  applies to signal  $z_j^{(0)}$  from neighbor  $j$ . Repeating the latter for  $k$  shifts, we get the recursion

$$\mathbf{z}^{(k)} = \Phi^{(k)} \mathbf{z}^{(k-1)} = \prod_{k'=0}^k \Phi^{(k')} \mathbf{x} = \Phi^{(k:0)} \mathbf{x}, \quad k = 0, \dots, K \quad (24)$$

where the product matrix  $\Phi^{(k:0)} = \prod_{k'=0}^k \Phi^{(k')} = \Phi^{(k)} \dots \Phi^{(0)}$  accounts for the weighted propagation of the graph signal  $\mathbf{z}^{(-1)} = \mathbf{x}$  from at most  $k$ -hops away neighbors. Each node is therefore free to adapt its weights for each iteration  $k$  to capture the necessary local detail.

**Edge varying filters and EdgeNets.** The collection of signals  $\mathbf{z}^{(k)}$  in (24) behaves like a sequence of parametric shifts, where at iteration  $k$  we use the parametric shift operator  $\Phi^{(k)}$  to shift-and-weigh the signal. Following the same idea as in (2), we can sum up *edge varying shifted* signals

$\mathbf{z}^{(k)}$  to get the input-output map  $\mathbf{y} = \mathbf{H}(\Phi)\mathbf{x}$  of an edge varying graph filter. For this relation to hold, the filter matrix  $\mathbf{H}(\Phi)$  should satisfy

$$\mathbf{H}(\Phi) = \sum_{k=0}^K \Phi^{(k:0)} = \sum_{k=0}^K \left( \prod_{k'=0}^k \Phi^{k'} \right). \quad (25)$$

The edge varying graph filter is characterized by the  $K + 1$  parameter matrices  $\Phi^{(0)}, \dots, \Phi^{(K)}$  and contains  $K(M + N) + N$  parameters. The edge varying graph filter forms the broadest family of graph filters: it generalizes the FIR filter in (2) (for  $\Phi^{(k:0)} = h_k \mathbf{S}^k$ ), the ARMA filter in (23), and almost all other filters employed to build GNNs including spectral filters [5], Chebyshev filters [6], Cayley filters, graph isomorphism filters, and also graph attention filters [21].

Substituting (25) into (8) yields an edge varying graph perceptron, which is the building block for edge varying GNNs or, for short, EdgeNets. EdgeNets are more than convolutional architectures; the high number of degrees of freedom and linear complexity render EdgeNets strong candidates for highly discriminatory GNNs in sparse graphs. If the graph is large, the EdgeNet can efficiently trade some edge detail (e.g., allowing edge varying weights only to a few nodes) to make the number of parameters independent of the graph dimension [16]. To control the number of parameters in EdgeNets we can adopt graph attention networks [21], see [16] for details on this and other alternatives.

## V. APPLICATIONS

We consider the application of GNNs for rating prediction in recommender systems (Section V-A) and learning decentralized controllers for flocking (Section V-B). These two applications aim at illustrating the use of GNNs in problems beyond semi-supervised learning.

We focus on the representation space of GNNs built with different filter *types* and compare them with that of linear FIR filters to corroborate the discussed insights. We note that, in all cases, the values of hyperparameters (number of layers  $L$ , filter taps  $K$  and features  $F_\ell$ ) are design choices that have been made after cross-validation<sup>2</sup>.

### A. Recommender Systems

Consider the problem of rating prediction in recommender systems. We have a database of users that have rated many items, and we use it to build a graph where each item is a node and

<sup>2</sup>The PyTorch GNN library used is available at <http://github.com/alelab-upenn/graph-neural-networks>.



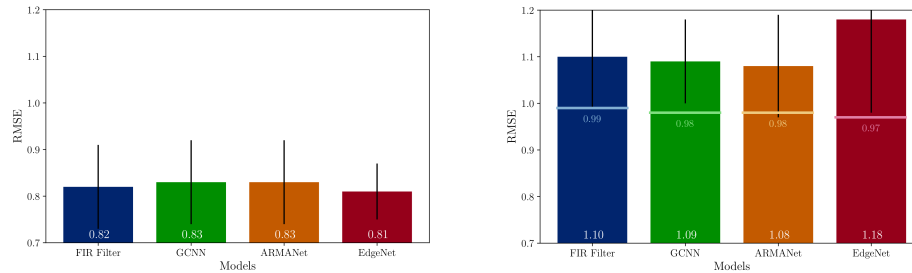


Figure 4. RMSE and standard deviation for different architectures in the movie recommendation problem. (Left) Training and testing on *Star Wars*. (Right) Training on *Star Wars*, testing on *Contact*; we also include the RMSE obtained for training and testing on the movie *Contact* as horizontal solid lines.

each edge weight is given by the rating similarity between items [22]. Then, given the ratings a specific user has given to some of the items, we want to predict the rating the same user would give to a specific item not yet rated. The ratings given by that user can be modeled as a graph signal, so that this becomes a problem of interpolating one of the (unknown) entries in it.

**Setup.** We consider items as movies and use a subset of the MovieLens-100k dataset, containing the 200 movies with the largest number of ratings [23]. The resulting dataset has 47,825 ratings given by 943 users to some of those 200 movies. The similarity between movies is the Pearson correlation [22, eq. (6)], which is further sparsified to keep only the ten edges with the stronger similarity. We split the dataset into 90% for training and 10% for testing. In this context, each user represents a graph signal, where the value at each node is the rating given to that movie. Movies not rated are given a value of zero. The objective is to estimate the rating a user would give to the movie *Star Wars* based on the ratings given by that same user to other movies and leveraging the graph of rating similarities.

**GNN models and training.** We implement a FIR graph filter (10)-(11) [cf. [22]], a GCNN (10)-(11)-(12), an ARMANet with  $T = 1$  Jacobi iterations (23), and an EdgeNet (25). The number of features in all cases is  $F_1 = 64$ , the filter order is  $K = 4$ , and ReLU nonlinearities are used. We include a local readout layer and extract the entry corresponding to the *Star Wars* movie as the estimate of the rating. The loss function is the smooth L1 loss  $J(x, y) = 0.5(x - y)^2$  if  $|x - y| < 1$  and  $|x - y| - 0.5$  otherwise, and the evaluation measure is the root MSE (RMSE). We train the architectures for 40 epochs with a batch size of 5 samples, using ADAM [18] with learning rate  $5 \times 10^{-3}$  and forgetting factors 0.9 and 0.999.

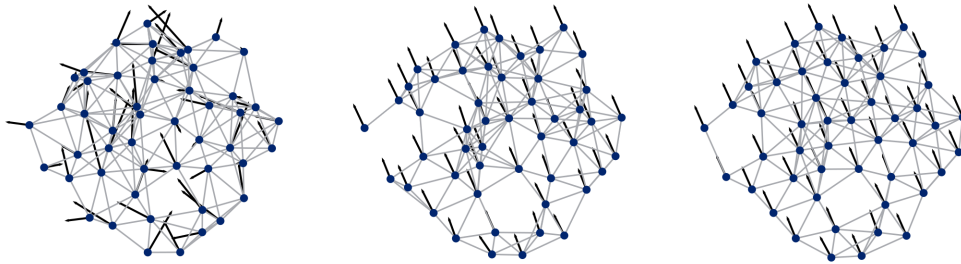


Figure 5. Snapshots of a sample trajectory. The dots illustrate the agents, the gray edges represent the communication links, and the arrows show the velocity. (Left) The agents start flying at time  $t = 0$ s with arbitrary velocities. (Middle) They manage to agree on a direction at  $t = 1$ s. (End) And they effectively fly together at  $t = 2$ s.

**Results.** Fig. 4 (left) shows that the rating for all models is similar, with the EdgeNet performing slightly better at an RMSE of  $0.81(\pm 0.05)$ . In Fig. 4 (right) we take the same models trained for estimating the rating for *Star Wars*, and extract the rating predicted for *Contact* instead. We do this in an attempt to show transferability of the trained models. In this case, the performance is similar among the FIR filter, GCNN and ARMANet, but the EdgeNet has severely degraded.

**Discussion.** First, we note that the GCNN performs similarly to the linear graph filter which, in light of Theorem 2, suggests that the relevant content is in low eigenvalues. Second, the GCNN and the ARMANet exhibit essentially the same performance, suggesting that the ARMA filter does not significantly increase the representation power, which is true in light of the FIR implementation (Jacobi) of ARMA filters. Third, the EdgeNet achieves the best performance when training and testing for the same movie, suggesting an increase in representation power, but does not transfer well to other settings, likely because it does not satisfy Theorems 1 nor 2.

### B. Learning Decentralized Controllers for Flocking

The objective of flocking is to coordinate a team of agents to fly together with the same velocity while avoiding collisions. Agents start flying at arbitrary velocities and need to take appropriate actions to flock together. This problem has a straightforward *centralized* solution that amounts to setting each agent’s velocity to the average velocity of the team [24, eq. (10)], but *decentralized* solutions are famously difficult to find [25]. Since GCNNs are naturally distributed, we use them to *learn* the decentralized controllers.

**Setup.** Let us consider a team of  $N = 50$  agents, in which each agent  $i$  is described at discrete time  $t$  by its position  $\mathbf{r}_i(t) \in \mathbb{R}^2$ , velocity  $\mathbf{v}_i(t) \in \mathbb{R}^2$ , and acceleration  $\mathbf{u}_i(t) \in \mathbb{R}^2$ . We want to control

Table I: Scalability. Trained on 50 agents. Tested on  $N$  agents. Optimal cost:  $51(\pm 1)$ .

$N$	50	62	75	87	100
FIR filter	408( $\pm 88$ )	408( $\pm 93$ )	434( $\pm 128$ )	420( $\pm 105$ )	430( $\pm 131$ )
GCNN	77( $\pm 3$ )	78( $\pm 3$ )	77( $\pm 2$ )	77( $\pm 2$ )	78( $\pm 2$ )

the acceleration  $\mathbf{u}_i(t)$  of each agent, so that they coordinate their velocities  $\mathbf{v}_i(t)$  to be the same for all  $i$ , see Fig. 5. We consider a decentralized setting where agents  $i$  and  $j$  can communicate with each other at time  $t$  only if  $\|\mathbf{r}_i(t) - \mathbf{r}_j(t)\| \leq R = 2\text{m}$ . This defines a communication graph  $\mathcal{G}(t)$  that changes with  $t$  as the agents move around, imposing a delayed information structure [24, eq. (2)]. Note that we can easily adapt graph filters (and thus, GCNNs) to the change in support matrices  $\mathbf{S}(t)$  by using delayed FIR filters  $\mathbf{H}(\mathbf{S}(t))\mathbf{x}(t) = \sum_{k=0}^K h_k \mathbf{S}(t)\mathbf{S}(t-1) \cdots \mathbf{S}(t-k+1)\mathbf{x}(t-k)$  [cf. (2)]. The filter taps  $h_k$  are the same, but the shift operators change with time. We generate 400 optimal trajectories for training and 20 for testing. We refer to [24] for details on system dynamics.

**GNN models and training.** We implement a linear FIR filter and a GCNN. We consider  $F_1 = 32$  features and filters of order  $K_1 = 3$ . We include a second, local readout layer to obtain the final acceleration  $\mathbf{u}_i(t) \in \mathbb{R}^2$  that each agent takes. We train the architectures using imitation learning by minimizing the MSE between the output action  $\mathbf{u}_i(t)$  and the optimal action  $\mathbf{u}_i^*(t)$  given by [24, eq. (10)]. At test time, we do not require access to the optimal action. The evaluation measure is the velocity variation of the team throughout the trajectory,  $N^{-1} \sum_t \sum_{i=1}^N \|\mathbf{v}_i(t) - \bar{\mathbf{v}}(t)\|^2$  with  $\bar{\mathbf{v}}(t) = N^{-1} \sum_{j=1}^N \mathbf{v}_j(t)$  being the average velocity at time  $t$ . We trained for 40 epochs with a batch size of 20 samples using ADAM [18] with learning rate  $5 \times 10^{-4}$  and forgetting factors 0.9 and 0.999.

**Results.** We observe in Table I the cost achieved by the GCNN-based controller is close to the optimal cost, while the FIR filter fails to control the system leading to a very high cost. We further investigate the effect of Theorems 1 and 2 by transferring at scale the learned solutions. That is, we take the controllers learned with teams of  $N = 50$  agents, and test them in teams of increasing size. The GCNN scales perfectly, maintaining the same performance.

**Discussions.** The GCNNs improved performance over the graph filter is expected since we know that optimal distributed controllers are nonlinear [25]. The GCNN also achieves a cost close to optimum, evidencing successful control. Once trained, this GCNN based controller can be used

in teams of arbitrary number of agents evidencing the properties of permutation equivariance and stability, and speaking to the potential of GCNNs for learning behaviors in homogenous teams.

## VI. CONCLUSION

Graph signal processing plays a crucial role in characterizing and understanding the representation space of graph neural networks. By emphasizing the role of graph filters and leveraging the concept of graph Fourier transform, we are able to derive fundamental properties such as permutation equivariance and stability, as well as establish a unified mathematical description. This reinforces the notion that GNNs are nonlinear extensions of graph filters, and thus GSP can help explain and understand the observed success of GNNs and contribute to improved designs.

As a matter of fact, several areas of interest lie ahead for GSP researchers to pursue. First, the understanding of what precise effect the nonlinearities have on the frequency content is limited. A better characterization of their effect in relation to the underlying topology is bound to help in designing appropriate ones. Second, the general relationship between the hyperparameters (number of layers, filter taps) and the characteristics of the graph (diameter, degree) is currently unknown. It is expected, for instance, the number of hops bears some relationship with the diameter of the graph, but no theoretical result is out there yet. Third, the bounds in the stability results are quite loose due to the coarse bound used on the eigenvectors. Thus, focusing on the eigenvector perturbation to improve the bound is a worthwhile pursuit. Fourth, the stability result holds for graphs of the same size. Extending this result to graphs of different size is an important research direction. Finally, we mention exploring the possibility of nonlinear aggregations of the filter banks, as well as using different shift operators at each layer.

From a higher vantage point, realizing GNNs are an object of study of GSP and regarding them as nonlinear extensions of graph filters, help us exploit our understanding of filtering techniques as well as leverage spectral domain analysis. Thus, GSP plays a crucial role in characterizing, understanding, and improving GNNs.

## REFERENCES

- [1] A. Ortega, P. Frossard, J. Kovačević, J. M. F. Moura, and P. Vandergheynst, "Graph signal processing: Overview, challenges and applications," *Proc. IEEE*, vol. 106, no. 5, pp. 808–828, May 2018.

- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, ser. The Adaptive Computation and Machine Learning Series. Cambridge, MA: The MIT Press, 2016.
- [3] S. Mallat, “Group invariant scattering,” *Commun. Pure, Appl. Math.*, vol. 65, no. 10, pp. 1331–1398, Oct. 2012.
- [4] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, “Geometric deep learning: Going beyond Euclidean data,” *IEEE Signal Process. Mag.*, vol. 34, no. 4, pp. 18–42, July 2017.
- [5] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, “Spectral networks and deep locally connected networks on graphs,” in *2nd Int. Conf. Learning Representations*, Banff, AB, 14–16 Apr. 2014, pp. 1–14.
- [6] M. Defferrard, X. Bresson, and P. Vandergheynst, “Convolutional neural networks on graphs with fast localized spectral filtering,” in *30th Conf. Neural Inform. Process. Syst.*, Barcelona, Spain, 5–10 Dec. 2016, pp. 3844–3858.
- [7] F. Gama, A. G. Marques, G. Leus, and A. Ribeiro, “Convolutional neural network architectures for signals supported on graphs,” *IEEE Trans. Signal Process.*, vol. 67, no. 4, pp. 1034–1049, Feb. 2019.
- [8] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” in *5th Int. Conf. Learning Representations*, Toulon, France, 24–26 Apr. 2017, pp. 1–14.
- [9] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, “Simplifying graph convolutional networks,” in *36th Int. Conf. Mach. Learning*, vol. 97. Long Beach, CA: PMLR, 9–15 June 2019, pp. 6861–6871.
- [10] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, “How powerful are graph neural networks?” in *7th Int. Conf. Learning Representations*, New Orleans, LA, 6–9 May 2019, pp. 1–17.
- [11] J. Atwood and D. Towsley, “Diffusion-convolutional neural networks,” in *30th Conf. Neural Inform. Process. Syst.*, Barcelona, Spain, 5–10 Dec. 2016.
- [12] F. Gama, J. Bruna, and A. Ribeiro, “Stability properties of graph neural networks,” *IEEE Trans. Signal Process.*, vol. 68, pp. 5680–5695, 25 Sep. 2020.
- [13] D. Zou and G. Lerman, “Graph convolutional neural networks via scattering,” *Appl. Comput. Harmonic Anal.*, vol. 49, no. 3, pp. 1046–1074, Nov. 2020.
- [14] E. Isufi, A. Loukas, A. Simonetto, and G. Leus, “Autoregressive moving average graph filtering,” *IEEE Trans. Signal Process.*, vol. 65, no. 2, pp. 274–288, Jan. 2017.
- [15] M. Coutino, E. Isufi, and G. Leus, “Advances in distributed graph filtering,” *IEEE Trans. Signal Process.*, vol. 67, no. 9, pp. 2320–2333, May 2019.
- [16] E. Isufi, F. Gama, and A. Ribeiro, “EdgeNets: Edge varying graph neural networks,” *arXiv:2001.07620v2 [cs.LG]*, 12 March 2020. [Online]. Available: <http://arxiv.org/abs/2001.07620>
- [17] M. Püschel and J. M. F. Moura, “Algebraic signal processing theory: Foundation and 1-d time,” *IEEE Trans. Signal Process.*, vol. 56, no. 8, pp. 3575–3585, Aug. 2008.
- [18] D. P. Kingma and J. L. Ba, “ADAM: A method for stochastic optimization,” in *3rd Int. Conf. Learning Representations*, San Diego, CA, 7–9 May 2015, pp. 1–15.
- [19] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*, vol. 323, no. 6088, pp. 533–536, Oct. 1986.
- [20] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, “CayleyNets: Graph convolutional neural networks with complex rational spectral filters,” *IEEE Trans. Signal Process.*, vol. 67, no. 1, pp. 97–109, Jan. 2019.

- [21] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, “Graph attention networks,” in *6th Int. Conf. Learning Representations*, Vancouver, BC, 30 Apr.-3 May 2018, pp. 1–12.
- [22] W. Huang, A. G. Marques, and A. Ribeiro, “Rating prediction via graph signal processing,” *IEEE Trans. Signal Process.*, vol. 66, no. 19, pp. 5066–5081, Oct. 2018.
- [23] F. M. Harper and J. A. Konstan, “The MovieLens datasets: History and context,” *ACM Trans. Interactive Intell. Syst.*, vol. 5, no. 4, pp. 19:(1–19), Jan. 2016.
- [24] E. Tolstaya, F. Gama, J. Paulos, G. Pappas, V. Kumar, and A. Ribeiro, “Learning decentralized controllers for robot swarms with graph neural networks,” in *Conf. Robot Learning 2019*, Osaka, Japan, 30 Oct.-1 Nov. 2019.
- [25] H. S. Witsenhausen, “A counterexample in stochastic optimum control,” *SIAM J. Control*, vol. 6, no. 1, pp. 131–147, 1968.