# FAST GRAPH CONVOLUTIONAL RECURRENT NEURAL NETWORKS

*Sai Kiran Kadambari and Sundeep Prabhakar Chepuri*

Indian Institute of Science, Bangalore, India
Email: {kadambarik;spchepuri}@iisc.ac.in

## ABSTRACT

**This paper proposes a Fast Graph Convolutional Neural Network (FGRNN) architecture to predict sequences with an underlying graph structure. The proposed architecture addresses the limitations of the standard recurrent neural network (RNN), namely, vanishing and exploding gradients, causing numerical instabilities during training. State-of-the-art architectures that combine gated RNN architectures, such as Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) with graph convolutions are known to improve the numerical stability during the training phase, but at the expense of the model size involving a large number of training parameters. FGRNN addresses this problem by adding a weighted residual connection with only two extra training parameters as compared to the standard RNN. Numerical experiments on the real 3D point cloud dataset corroborates the proposed architecture.**

***Index Terms***— **Deep learning, Graph neural networks, graph signal processing, recurrent neural networks.**

## I. INTRODUCTION

Many real-world datasets occur as structured sequences, e.g., space-time series, or graph/network-time series. Typical examples of such data are multichannel speech data, videos, dynamic point clouds, timeseries data from brain networks or social networks, where the successive frames might have a dynamic pattern and each time frame might have a spatial or graph structure. Developing generative models to process such structured time-varying data for prediction or interpolation is of significant interest in data analytics.

Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) are two popular variants of neural networks commonly used in a wide variety of engineering and science applications. CNN comprises of a convolutional layer and pooling layer as hidden layers and is known to capture intricate structures in the data. Whereas, RNNs have a memory element that captures the information about the past and decisions are based on the gathered memory.

Architectures combining RNN and CNN have been proposed to find patterns in time-varying data [1]–[3]. These

models assume that the data is defined on a regular domain, and the convolutions in this case are simple 2-D convolutions. In many cases, the data available might not be supported on a regular domain [4], [5]. Examples of data residing on irregular domains are data from weather monitoring stations, biological networks, and social networks, to list a few. Graphs may be used to represent the irregular domain on which the data is defined and explain the complex relationships in such data [6], [7]. More specifically, the data is indexed by the nodes of the graph, and the edges encode the pairwise relationships between the nodes. Numerous approaches have been proposed to learn the underlying graph structure from the available data [8], [9].

The standard RNN architecture is known to have vanishing or exploding gradients, which makes it unstable during training. A particular class of RNN that overcomes this problem is the Long Short-Term Memory (LSTM) architecture [10]. These variants of RNNs have a gated architecture with a large number of training parameters. As the number of trainable parameters in LSTM is large, they take more training time, even on high-end computational machines, due to the required memory.

Graph Convolutional Neural networks (GCNN) that generalize the CNNs to handle graph-structured data on arbitrary graphs have been proposed in [11], [12], where the convolution operator is now generalized using polynomials of the graph Laplacian matrix. When the graph data is time-varying (e.g., dynamic 3D point cloud data), to capture the temporal variations of such data, LSTM combined with GCNN is proposed in [13], where GCNNs capture the spatial structure and LSTMs capture the temporal variations of the data.

In this paper, we propose a stable architecture of RNN that combines GCNN with a standard RNN. To stabilize the gradients a weighted residual connection is introduced. The proposed architecture has a much fewer number of training parameters as compared to architectures that combine LSTM with GCNN. Hence the term fast in FGRNN. The proposed architecture is inspired by [14], where they propose a stable, scalable, and a faster variant of RNN for data defined on regular domains. The main contributions of this paper are as follows. We propose a stable FGRNN architecture for efficient training and prediction of the data defined on irregular domains. We show that the proposed FGRNN

architecture with only 2 additional parameters as compared to the standard RNN is stable during the training phase and overcomes the vanishing or exploding gradients problem. The experiments on real 3D point cloud data, where the task is to predict the next 3D point cloud frame, demonstrate the developed theory.

Throughout this paper, we will use upper and lower case boldface letters to denote matrices and column vectors, respectively. $\text{diag}[\cdot]$ is a diagonal matrix with its argument along the main diagonal. $X_{ij}$ and $x_i$ denote the $(i,j)$th element and $i$th element of $\mathbf{X}$ and $\mathbf{x}$, respectively.

## II. PROBLEM STATEMENT

Consider a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $N$ vertices (nodes), where $\mathcal{V}$ is the set of vertices (or nodes) and $\mathcal{E}$ denotes the edge set such that $(i,j) \in \mathcal{E}$ if there is a connection between node $i$ and node $j$. The structure of the graph with $N$ nodes is captured by the weighted adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ whose $(i,j)$th entry denotes the weight of the edge between node $i$ and node $j$. This matrix represents the network connectivity. We assume that the graph is undirected with positive edge weights. The corresponding graph Laplacian matrix is a symmetric matrix of size $N$, given by $\mathbf{L} = \mathbf{I}_N - \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$, where $\mathbf{D} \in \mathbb{R}^{N \times N}$ is the diagonal matrix whose diagonal entries are given by the length $N$ degree vector $\mathbf{d} = \mathbf{A1}$. We call the set of graph signals $\{\mathbf{x}_t\}_{t=1}^{T}$, indexed by the vertices of the graph $\mathcal{G}$, collected in the $N \times T$ matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_T]$ as graph data. That is, $\mathbf{x}_t \in \mathbb{R}^N$ is the graph signal at time $t$.

The goal of this paper is to learn a function for sequence modeling, where the task is to predict the most likely feature vector based on the previous observations. More percisely, given the previous $T$ observations of the data, we are interested in predicting the most likely feature vector $\mathbf{x}_{t+1}$:

$$\hat{\mathbf{x}}_{t+1} = \underset{\mathbf{x}_{t+1}}{\arg\max} \quad P(\mathbf{x}_{t+1}|\mathbf{x}_{t-T}, ..., \mathbf{x}_t), \tag{1}$$

where in (1) we maximize the likelihood of $\mathbf{x}_{t+1}$ given the previous $T$ observations. Such problems usually appear in language modeling, image prediction, to name a few, wherein the task is to predict the most likely feature vector given the previous $T$ observations. In this paper, we focus on the case where the data at each time instance has a structure, which is defined by the graph $\mathcal{G}$.

## III. GRAPH CONVOLUTIONS

As the graphs do not have a natural ordering, the standard convolution operation cannot be generalized to arbitrary graphs using localized filters. Hence, a spectral definition of graph convolution is defined in [15] and is based on an elementwise multiplication in the graph frequency domain. Using this definition, the graph convolution of the graph signal $\mathbf{x}$ is given as $\mathbf{y} = g_\theta \star_\mathcal{G} \mathbf{x} = \mathbf{U}g_\theta(\mathbf{\Lambda})\mathbf{U}\mathbf{x}$, where $\mathbf{U} \in \mathbb{R}^{N \times N}$ is the matrix of eigenvectors, $\mathbf{\Lambda} \in \mathbb{R}^{N \times N}$ is the matrix of eigenvalues of the graph Laplacian matrix

$\mathbf{L}$. Here, $g_\theta(\mathbf{\Lambda})$ represents the responce of the graph filter in the frequency domain and $\star_\mathcal{G}$ is the graph convolution operator. This method is computationally expensive, as it involves multiplications with a dense eigenvector matrix $\mathbf{U}$. To circumvent this issue, [11] defined the graph convolution operation by approximating $g_\theta(\mathbf{\Lambda})$ with a truncated Chebyshev polynomial expansion of order $K$. Mathematically, the graph convolution operation is given by

$$g_\theta \star_\mathcal{G} \mathbf{x} = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\mathbf{L}})\mathbf{x}, \tag{2}$$

where the parameter $\boldsymbol{\theta} = [\theta_0, \theta_1, ...\theta_{K-1}]^\mathrm{T} \in \mathbb{R}^K$ is a vector of Chebyshev coefficients and $T_k(\tilde{\mathbf{L}}) \in \mathbb{R}^{N \times N}$ is the Chebyshev polynomial of order $K$ evaluated at the normalized Laplacian matrix defined as $\tilde{\mathbf{L}} = 2\mathbf{L}/\lambda_{max} - \mathbf{I}_N$. The above graph convolution operation incurs linear complexity, i.e., $\mathcal{O}(K|\mathcal{E}|)$, due to the recurrence relation $T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$. Here, $T_0 = 1$ and $T_1 = x$.

A first order approximation $(K = 1)$ of the (2) is proposed in [12]. Using the approximation in [12], the convolution operation simplifies to

$$g_\theta \star_\mathcal{G} \mathbf{x} \quad = \mathbf{W}\tilde{\mathbf{L}}_1\mathbf{x}, \tag{3}$$

where $\tilde{\mathbf{L}}_1 = \mathbf{I}_N + \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$, $\mathbf{W} \in \mathbb{R}^{P \times N}$ is the filter parameter and $P$ is the hidden state dimension. While we focus on the graph convolution framework introduced in [12] for the analysis of the proposed architecture, we demonstrate the effectiveness of the proposed model using both (2) and (3).

## IV. RECURRENT NEURAL NETWORKS

For tasks like sequence modeling, standard RNNs are typically used. The standard RNN maintains a hidden state that captures the temporal variations in the data. The hidden state in the standard RNN at the time instance $t$ is given by $\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b})$, where $\mathbf{W} \in \mathbb{R}^{P \times N}$, $\mathbf{U} \in \mathbb{R}^{P \times N}$, and $\mathbf{b} \in \mathbb{R}^{P \times 1}$ are the training parameters. Here, $\sigma(\cdot)$ is the nonlinear activation function, and typical choices for $\sigma(\cdot)$ are such as $\tanh$, ReLU, sigmoid.

It is well known that standard RNNs suffers from exploding and vanishing gradient issues, due to which they are highly unstable during the training phase. A residual connection is introduced in [14] to alleviate this numerical instability. This fast variant of RNN in short, (FRNN) has much fewer parameters as compared to LSTM. Mathematically, the hidden states of FRNN is given by

$$\begin{aligned} \tilde{\mathbf{h}}_t &= \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}) \\ \mathbf{h}_t &= \alpha\tilde{\mathbf{h}}_t + \beta\mathbf{h}_{t-1}, \end{aligned} \tag{4}$$

where $\alpha$ and $\beta$ are scalar the training parameters. FRNN updates the hidden states in a controlled manner using only 2 extra training parameters $\alpha$ and $\beta$. Moreover, the number of additional computations required per time step is $N$, which is usually a tiny fraction as compared to the operations in RNN

and as such it is very small compared to the computational complexities of the gated architectures like LSTM and GRU.

## V. FAST GRAPH RECURRENT NEURAL NETWORKS

In this section, we propose a Fast Graph Convolutional Recurrent Neural Network (FGRNN) architecture that combines FRNN with GCNNfor the task of prediction. For the task of sequence modeling for data defined on regular domains, architectures that combine RNN and CNN are proposed in [16]. In this architecture, LSTM is used, where the 2D convolutions replace the multiplications with the dense matrices. Inspired by this approach, when the data is defined on irregular domains, we propose to replace the multiplications by the dense matrices in FRNN with graph convolutions. Specifically, the update equations for the hidden state is given by

$$
\begin{aligned}
\tilde{\mathbf{h}}_t &= \sigma(\mathbf{W} \star_{\mathcal{G}} \mathbf{x}_t + \mathbf{U} \star_{\mathcal{G}} \mathbf{h}_{t-1} + \mathbf{b}) \\
\mathbf{h}_t &= \alpha \tilde{\mathbf{h}}_t + \beta \mathbf{h}_{t-1},
\end{aligned} \tag{5}
$$

where recall that $\star_{\mathcal{G}}$ is the graph convolution operator defined in (2) or (3). We focus on (3), for the sake of simplicity, for which the hidden state at time $t$ smiplifies to

$$
\begin{aligned}
\tilde{\mathbf{h}}_t &= \sigma(\mathbf{W}\mathbf{L}\mathbf{x}_t + \mathbf{U}\mathbf{L}\mathbf{h}_{t-1} + \mathbf{b}), \\
\mathbf{h}_t &= \alpha \tilde{\mathbf{h}}_t + \beta \mathbf{h}_{t-1}.
\end{aligned} \tag{6}
$$

Given the hidden state at any time step $t$, we predict the next feature using the linear relation given by $\hat{\mathbf{x}}_{t+1} = \mathbf{V}\mathbf{h}_t + \mathbf{z}$, where $\mathbf{V} \in \mathbb{R}^{N \times P}$ and $\mathbf{z} \in \mathbb{R}^N$ are the training parameters. At each time step, we define the prediction loss function as $J_t(\hat{\mathbf{x}}_{t+1}, \mathbf{x}_{t+1}, \boldsymbol{\Theta}) = \|\mathbf{x}_{t+1} - \hat{\mathbf{x}}_{t+1}\|_2^2$, where $\boldsymbol{\Theta} = \{\mathbf{W}, \mathbf{U}, \mathbf{V}, \alpha, \beta, \mathbf{b}, \mathbf{z}\}$ is the set of all training parameters. The prediction loss function after $T$ time steps is given by $J_T = \sum_{t=1}^{T} J_t(\mathbf{x}_{t+1}, \hat{\mathbf{x}}_{t+1}, \boldsymbol{\Theta})$.

To update the training parameters during the training phase using backpropagation through time (BPTT), we calculate the gradient of the loss function with respect to training parameters. For simplicity, assume that $\boldsymbol{\Theta} = \{\mathbf{W}, \mathbf{U}, \mathbf{V}\}$ are the trainig parameters. Thus, the gradients of the loss functions $J_T$ w.r.t. parameters $\mathbf{W}$, $\mathbf{U}$, and $\mathbf{V}$ are given by

$$
\begin{aligned}
\frac{\partial J_T}{\partial \mathbf{W}} &= \sum_{t=1}^{T} \frac{\partial J_t}{\partial \mathbf{h}_T} \prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}, \\
\frac{\partial J_T}{\partial \mathbf{U}} &= \sum_{t=1}^{T} \frac{\partial J_t}{\partial \mathbf{h}_T} \prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_1}{\partial \mathbf{U}}, \\
\frac{\partial J_T}{\partial \mathbf{V}} &= \sum_{t=1}^{T} \frac{\partial J_t}{\partial \mathbf{h}_T} \prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_1}{\partial \mathbf{V}},
\end{aligned} \tag{7}
$$

where the term $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \alpha \mathbf{D}_t \mathbf{U} \mathbf{L} + \beta \mathbf{I}_N \in \mathbb{R}^{N \times N}$ common to all the gradients determines the numerical stability during the training phase. This term is multiplied by itself $T - 2$ times,and depending on the conditioning of $\mathbf{U}$ and/or

$\mathbf{L}$, $\prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$ may quickly become ill-conditioned with $\beta = 0$. Here, $\mathbf{D}_t$ is the diagonal matrix with pointwise nonlinearity given by $\mathbf{D}_t = \text{diag}(\sigma'(\mathbf{W}\mathbf{L}\mathbf{x}_t + \mathbf{U}\mathbf{L}\mathbf{h}_t + \mathbf{b}))$ with $\sigma'(\cdot)$ being the gradient of $\sigma(\cdot)$. As the special case, if the activation function is ReLU, then the matrix $\mathbf{D}_t$ is an identity matrix. Notice that when $\alpha = 1$ and $\beta = 0$, we have the standard RNN architecture.

To analyze the problem of the vanishing and exploding gradients in the standard RNN combined with the graph convolution operation, let us assume that all the training parameters are scalars and the activation function is ReLU. Then the term $\prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$ simplifies to

$$
\prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = (u\mathbf{L})^{T-2}. \tag{8}
$$

The stability of the gradient depends on the largest eigenvalue of (8). If its value is sufficiently small, i.e., $< 1$ the gradient will shrink exponentially. Moreover, if its value is large, the gradient will explode. For $\tilde{\mathbf{L}}$, the eigenvalues in the range $[0, 2]$, thus gradient for the standard RNN, will be an exponential in $T$. This implies that, relative to the largest eigenvalue, the gradient may explode or vanish exponentially, leading to numerical instabilities during training.

In order to stabilize the gradients, we add a weighted residual connection, with two trainable parameters $\alpha$ and $\beta$ as in (6) for which (8) now becomes $\prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = (\alpha \mathbf{U} \mathbf{L} + \beta \mathbf{I}_N)^{T-2}$. By appropriately choosing $\alpha$ and $\beta$, the gradients may be stabilized. More generally, the condition number of $\prod_{t=2}^{T} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$, which is bounded as

$$
M \le \frac{(1 + \frac{\alpha}{\beta} \max_t \|\mathbf{D}_t \mathbf{U} \mathbf{L}\|_F^2)^{T-2}}{(1 - \frac{\alpha}{\beta} \max_t \|\mathbf{D}_t \mathbf{U}.\mathbf{L}\|_F^2)^{T-2}} \tag{9}
$$

This determines the stability of the gradient. In contrast to the standard RNN, if $\beta = 1$ and $\alpha = 0$, then the condition number of $M$ is bounded 1. Thus leading to stable gradients during training, but completely ignores the training data. So, we allow the parameters $\alpha$ and $\beta$ to be trained such that FGRNN updates the hidden states in a controlled manner. In other words, the parameters $\alpha$ and $\beta$ limit the extent to which the current feature vector $\mathbf{x}_t$ updates the hidden states $\mathbf{h}_t$. For instance, if $\beta \approx T \max_t \|\mathbf{D}_t \mathbf{U} \mathbf{L}\|_F^2$ and $\alpha \approx 1 - \beta$, then $M = \mathcal{O}(1)$. Also, the FGRNN controls the condition number of the gradient using only two additional parameters as compared to the standard RNN. Furthermore, the existing unitary RNN methods for the data defined on a regular domain are motivated by a similar observation, where they control the stability of the gradients by restricting $\mathbf{U}$ to a unitary matrix. However, while dealing with graph data, $\mathbf{D}_t \mathbf{U} \mathbf{L}$ might still become ill-conditioned even with an unitary $\mathbf{U}$. Thus they might still have the vanishing gradient problem. The proposed method allows the residual weights $\alpha$ and $\beta$ to be trained such that the condition number $M$ is

**Fig. 1**: *Dynamic 3D point cloud dataset*. The colored dot indicates a node of the 3D point cloud data. Each image shows the predicted 3D point cloud frame using different architectures. (a). Ground truth (b). Proposed model (filter based on [11]) (c). Proposed model (filter based on [12]) (d). LSTM with graph regularizer (e). LSTM with GCN [13].



**Fig. 2**: *Test loss of the 3D point cloud dataset*: Left: Test loss of the various architectures trained on the 3D point cloud dataset. Right: Variations of the training parameters $\alpha$ and $\beta$ with time $T$.

restricted and thus prevents the numerical instabilities due to vanishing and exploding gradients.

## VI. NUMERICAL ANALYSIS

In this section, we test the proposed FGRNN architecture and compare its performance with the traditional architectures that combine GCNs with a stable variant of RNN, namely, LSTM. We compare the proposed method with (1). the baseline method that combines GCN with LSTM [13], (2). LSTM with a graph regularizer in the loss function, (3). proposed method, i.e., FGRNN with different convolution operators, i.e., Chebyshev polynomial (2) and its first order approximation (3).

We use the dynamic 3D point cloud dataset of a human pose. Each human pose is captured using $1502$ 3D data points and there are $573$ such frames. This corresponds to a graph data, where the data points correspond to the 3D displacement of the nodes of the underlying graph. Each data frame at time instance $t$ is given by $\mathbf{X}_t \in \mathbb{R}^{N \times 3}$ where $N = 1502$ is the number of nodes. We use $80\%$ of the frames for training the network and the remaining $20\%$ of the time frames to test the performance.

For this dataset, we are interested in predicting the most likely next 3D point frame (here the next human pose)

given the previous $T$ frames of data as in (1), where $P(\mathbf{X}_{t+1}|\mathbf{X}_{t-T}, ..., \mathbf{X}_t)$ models the probability of the frame $\mathbf{X}_{t+1}$ given the past $T$ observations. As the task we are interested in is the prediction, we define the loss occured for predicting the next data frame at time $t$ as $\|\mathbf{X}_{t+1} - \hat{\mathbf{X}}_{t+1}\|_F^2$, where each $\mathbf{X}_{t+1}$ is a 3D data point. All the architectures (except LSTM with a graph regularizer) are trained by minimizing the aforementioned loss function using BPTT. For LSTM with a graph regularizer architecture, we define the loss function as $\|\mathbf{X}_{t+1} - \hat{\mathbf{X}}_{t+1}\|_F^2 + \lambda \mathrm{tr}(\mathbf{X}_{t+1}^T \mathbf{L} \mathbf{X}_{t+1})$, where $\lambda$ is a positive regularizer chosen based on a grid search that leads to the best test loss. We construct the Laplacian matrix $\mathbf{L}$ from the training data using K-nearest-neighbor such that each node has a degree of $6$. All the architectures are trained using ADAM optimizer with a learning rate of $10^{-2}$ and a decay rate of $0.9$. All the models are implemented in Tensorflow: r1.15 and each model is run for 10 epochs and their test loss is shown in Fig. 2a.

For a fair comparison with [16], all the experimental models are implemented with the model defined in (5), which is the same as replacing the 2D convolutions by graph convolutions. Fig. 2 shows the performance of the various models implemented, and we observe that all the models converge before 10 epochs. These results show the ability

| | # parameters | point cloud |
|---|---|---|
| Standard FRNN | $3N^2 + 2N + 2$ | 6771018 |
| Proposed (filter based on [11]) | $3K + 2N + 2$ | **3015** |
| Proposed (filter based on [12]) | $3P^2 + 2N + 2$ | **3033** |
| LSTM without GCN | $8N^2 + 4N$ | 18054040 |
| LSTM with GCN  [13] | $4N + 8K$ | 6032 |

**Table I**: Comparison between the models in terms of number of trainable parameters.

of the proposed method to capture the structure in the graph time-series. We can observe from Fig. 2a that FGCNNs implemented with graph filter based on [12] and [11] offer better performance than regular LSTMs with a graph regularizer and architecture that combines GCN with LSTM [16]. In Fig. 1, we illustrate the 3D point cloud frames of the ground truth, and the estimated frames from various architectures we have considered. We see that the predicted 3D point cloud frame by the proposed method is more consistent with the ground truth than the frame estimated by LSTM with a graph regularizer. Moreover, the predicted feature is visually the same as the baseline methods (namely, LSTMs with a graph convolutions). This demonstrates that the problem of vanishing and exploding gradients can be overcome by the addition of a simple weighted residual connection to the standard RNN, which means that the proposed FGRNN is stable and can be trained efficiently.

Table I shows the computational complexity in terms of the number of training parameters for the different considered methods. We can see that the proposed method is computationally efficient than any other baseline methods as the number of trainable parameters in the proposed architecture is the least. This demonstrates that the proposed FGRNN models are accurate and faster to train as compared to the baseline models. Finally, Fig. 2b shows the learnt $\alpha$ and $\beta$ on the datasets with $T$ time steps. It is clear from the figure that the learned $\beta$ is a decreasing function of $T$. Moreover, $\alpha$ can be seen close to $1 - \beta$ for large $T$, while corroborates the FGRNNs theoretical analysis.

## VII.  CONCLUSIONS

This paper proposes a Fast Graph Recurrent Neural Network (FGRNN) architecture for efficient training and prediction of data defined on irregular (non-euclidean) domains. The sandard RNN architecture is known to be unstable during training due to vanishing and exploding gradients. Hence, gated architectures, namely, LSTMs and GRUs, are proposed to alleviate this issue at the cost of the computational complexity. FGRNN architecture is obtained by incorporating a weighted residual connection with only two scalar parameters into the standard RNN architecture. FGRNN model has a fewer number of training parameters, lesser training times, and is more stable than the standard RNN. These architectures can match the state-of-the-art

gated RNN architectures with a significantly lower number of training parameters and computational cost.

## VIII. REFERENCES

[1] J. Donahue, L. A. Hendricks, M. Rohrbach, S. Venugopalan, S. Guadarrama, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 677–691, Apr. 2017.

[2] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 664–676, Apr. 2017.

[3] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2015.

[4] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, "Geometric deep learning: Going beyond euclidean data," *IEEE Signal Processing Magazine*, vol. 34, no. 4, pp. 18–42, July 2017.

[5] F. Monti, M. Bronstein, and X. Bresson, "Geometric matrix completion with recurrent multi-graph neural networks," in *Advances in Neural Information Processing Systems*, Sep. 2017, pp. 3697–3707.

[6] A. Sandryhaila and J. M. F. Moura, "Big data analysis with signal processing on graphs: Representation and processing of massive data sets with irregular structure," *IEEE Signal Processing Magazine*, vol. 31, no. 5, pp. 80–90, Sep. 2014.

[7] A. Sandryhaila and J. M. Moura, "Discrete signal processing on graphs: Frequency analysis." *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3042–3054, 2014.

[8] X. Dong, D. Thanou, M. Rabbat, and P. Frossard, "Learning graphs from data: A signal representation perspective," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 44–63, May 2019.

[9] S. P. Chepuri, S. Liu, G. Leus, and A. O. Hero III, "Learning sparse graphs under smoothness prior," in *Proc. of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, New Orleans, USA, Mar. 2017.

[10] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[11] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Advances in neural information processing systems*, 2016.

[12] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations (ICLR)*, 2017.

[13] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson, "Structured sequence modeling with graph convolutional recurrent networks," *Lecture Notes in Computer Science*, p. 362373, 2018.

[14] A. Kusupati, M. Singh, K. Bhatia, A. Kumar, P. Jain, and M. Varma, "Fastgrnn: A fast, accurate, stable and tiny kilobyte sized gated recurrent neural network," in *Advances in Neural Information Processing Systems*, 2018.

[15] J. Bruna, W. Zaremba, A. Szlam, and Y. Lecun, "Spectral networks and locally connected networks on graphs," Dec. 2013.

[16] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," in *Advances in neural information processing systems*, 2015.