
A Multiscale Graph Convolutional Network Using Hierarchical Clustering

Alex Lipov¹ Pietro Lió¹

Abstract

The information contained in hierarchical topology, intrinsic to many networks, is currently underutilised. A novel architecture is explored which exploits this information through a multiscale decomposition. A dendrogram is produced by a Girvan-Newman hierarchical clustering algorithm. It is segmented and fed through graph convolutional layers, allowing the architecture to learn multiple scale latent space representations of the network, from fine to coarse grained. The architecture is tested on a benchmark citation network, demonstrating competitive performance. Given the abundance of hierarchical networks, possible applications include quantum molecular property prediction, protein interface prediction and multiscale computational substrates for partial differential equations.

1. Introduction and Related Work

Most networks display some level of hierarchical topology (Barabási & Pósfai, 2016); examples include actor networks, the semantic web and the internet at the autonomous system level (Ravasz & Barabási, 2003).

The defining feature of hierarchical network topology is how the clustering coefficient, C_i , of the i -th node scales with its degree, k_i ; namely,

$$C_i(k) \propto k_i^{-1}. \quad (1)$$

Nodes with high degree and low clustering coefficient are those that connect different communities. Nodes with low degree and high clustering coefficient are those that connect within the community more than to different communities. Another distinct characteristic of hierarchical networks is that the clustering coefficient is independent of the number

¹Department of Computer Science and Technology, University of Cambridge, UK. Correspondence to: Alex Lipov <al822@cam.ac.uk>.

of nodes and that they show scale-free topology (Barabási & Pósfai, 2016).

In this work, we present and explore a novel architecture operating on the graph domain which attempts to utilise the intrinsic hierarchical topological information that is embedded in a wide range of datasets. Graph neural networks, a booming field of deep learning, was built on the assertion that intrinsic graph structure in datasets was underutilised (Zhou et al., 2018). However, the field of multiscale graph neural networks effectively does not exist, despite the fact that the evidence supporting the claim of hierarchical underutilisation is just as strong as for graph structure underutilisation (Ravasz & Barabási, 2003).

Limited work has been done in the area of multiscale graph neural networks. Luzhnica, Day & Liò (2019) introduced clique pooling; however cliques (complete subgraphs) are more restrictive than clusters (dense subgraphs) — it is more likely that hierarchical networks consist of clusters rather than cliques (beyond triangular cliques). The work on multiscale graph convolutions by Haija *et al.* (2018) is the closest piece of work to ours; however it does not use hierarchical clustering algorithms.

In terms of applications; Lu *et al.* (2019) used a hierarchical level-by-level treatment of quantum interactions in molecules, Zitnik *et al.* (2017) used hierarchical multiplex graphs for multiple spatial scales of brain tissue and Kim *et al.* (2019) developed a temporal multiscale graph convolutional neural network for analysing bike-sharing demands. Our multiscale decomposition was loosely inspired by the multiscale neural network using hierarchical block matrices, by Fan *et al.* (2019), however it was designed to solve partial differential equations and does not operate on the graph domain.

2. Architecture

The architecture is shown in Figure 1. We input an unweighted, undirected, graph, $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where there are N nodes, $\mathcal{V} \in \{\mathcal{V}_1, \dots, \mathcal{V}_N\}$, and M edges, $\mathcal{E} \in \{\mathcal{E}_1, \dots, \mathcal{E}_M\}$. We construct $\hat{A} = A + I$, where $I \in \mathbb{R}^{N \times M}$ is the identity matrix and $A \in \mathbb{R}^{N \times N}$ is the adjacency ma-

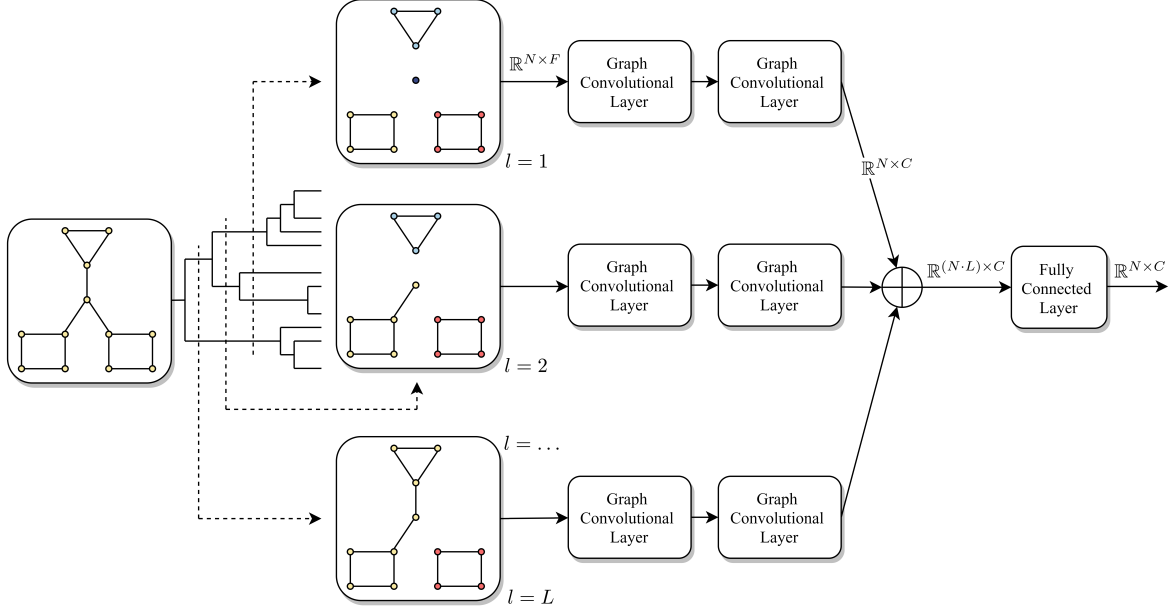


Figure 1. Our architecture which exploits hierarchical information through a multiscale decomposition. A dendrogram is produced by a Girvan-Newman hierarchical clustering algorithm, as shown on the left. It is segmented and fed through graph convolutional layers, shown in the middle, allowing the architecture to learn multiple scale latent space representations of the network, from fine to coarse grained. These representations are then combined, shown on the right, and fed through a fully connected layer which outputs the model’s node classification predictions.

trix given by,

$$A_{ij} = A_{ji} = \begin{cases} 1 & \text{if node } i \text{ links node } j \\ 0 & \text{otherwise} \end{cases}. \quad (2)$$

A multiscale decomposition, f , is taken through hierarchical clustering,

$$f(\mathcal{G}) = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_L\}, \quad (3)$$

and hence,

$$f(\hat{A}) = \{\hat{A}_l\} \text{ where } l \in \{1, \dots, L\}. \quad (4)$$

We pass this and the input feature matrix, $X \in \mathbb{R}^{N \times F}$, with F features per node, through the first convolutional layer, with F input nodes, obtaining the set over l ,

$$\left\{ H_l^{(1)} : H_l^{(1)} = \sigma \left(\hat{D}_l^{-\frac{1}{2}} \hat{A}_l \hat{D}_l^{-\frac{1}{2}} X W_l^{(0)} \right) \right\}, \quad (5)$$

where $H_l^{(1)} \in \mathbb{R}^{N \times F'}$, F' is the number of input nodes to the second layer, $W_l^{(q)}$ is the weight matrix for the q -th neural network layer on the l -th scale, σ is a non-linear activation function and \hat{D}_l is the l -th diagonal node degree matrix. The Kipf & Welling (2016) propagation rule was used in equation 5. We then pass through the second convolutional layer,

$$\left\{ H_l^{(2)} : H_l^{(2)} = \sigma \left(\hat{D}_l^{-\frac{1}{2}} \hat{A}_l \hat{D}_l^{-\frac{1}{2}} H_l^{(1)} W_l^{(1)} \right) \right\}, \quad (6)$$

where $H_l^{(2)} \in \mathbb{R}^{N \times C}$, with C being the number of classes. We then vertically concatenate the latent feature matrices from each scale, obtaining

$$H^{(3)} = H_1^{(2)} \oplus H_2^{(2)} \oplus \dots \oplus H_L^{(2)}. \quad (7)$$

Here, $H^{(3)} \in \mathbb{R}^{(N-L) \times C}$. Finally this is fed through a fully connected (FC) layer, resulting in the node classification probability matrix as output,

$$Y = \sigma(WH^{(3)} + b), \quad (8)$$

where $Y \in \mathbb{R}^{N \times C}$, the weight matrix $W \in \mathbb{R}^{N \times (N-L)}$ and b is a bias matrix $Y \in \mathbb{R}^{N \times C}$. In our implementation of this architecture, we chose $F' = (F + C)/2$.

We used a graph convolutional network (GCN) layer depth of two here but this architecture can easily be generalised to any GCN layer depth. Ensembling, as we have done here, has significant literature demonstrating its power as an architectural choice (Xibin et al., 2020).

3. Experiments

Tests were conducted on a truncated version of the Cora dataset (<https://relational.fit.cvut.cz/dataset/CORA>). Its properties are detailed in Table 1. Code can be found at [link removed for review].

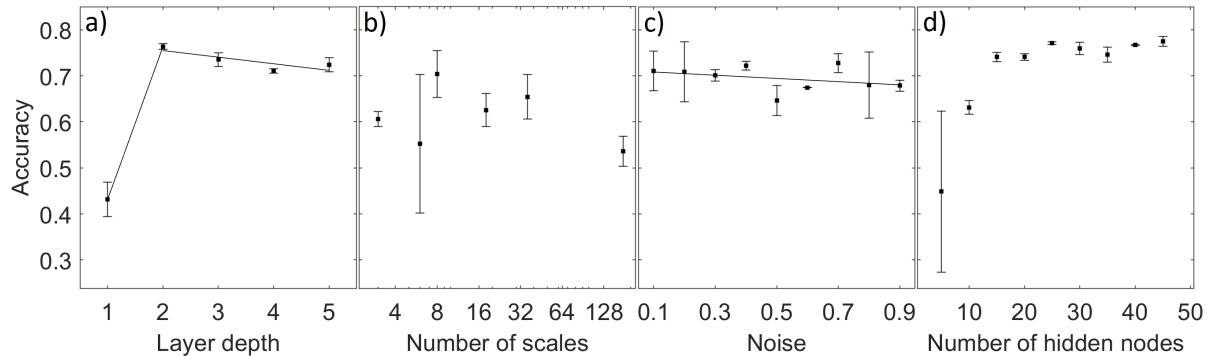


Figure 2. Node classification accuracy of the architecture when the GCN layer depth, number of decomposition scales and number of hidden nodes is varied. Accuracy upon corrupting the training set with random binary valued feature vectors is also shown, where the noise is altered. There is an optimal depth, shown by (a). Accuracy peaks at a threshold scale shown by (b), but this trend is somewhat weak. Robustness to noise is demonstrated by (c) and a threshold neural network width is required from observing (d).

Table 1. Key properties of the truncated Cora dataset, used in the experiments.

Truncated Cora Dataset	
Type	Machine learning citation network
Task	Semi-supervised node classification
Number of nodes	2708
Number of labelled nodes	140
Number of edges	5429
Number of classes	7
Feature vector length	1433

3.1. Hyperparameters

Prior to performing the following experiments, the Cora network was put through a Girvan-Newman clustering algorithm. It produced 528 usable networks (excluding the last one which contained no edges, hence graph convolutions could not be performed properly). The non-linear activation function used was ReLU. Adam optimization was used with a learning rate of 0.01 and weight decay of 5×10^{-4} . The number of epochs used to train the convolutional layers was 300 for each graph. The number of epochs used to train the FC layer was 10. Given that this was a multi-class classification task with no major skew to the dataset, accuracy was deemed to be a sufficient performance metric. Two tests were performed for each data point, with error bars given by the standard deviation. The strategy for assigning the number of hidden nodes for a given GCN layer was to take the arithmetic mean of the number of nodes in the directly neighbouring layers, ensuring consistency. Each data point took approximately 10 minutes on a Nvidia GTX Titan X. Figure 2 summarises the results. Overall, the architecture achieved accuracies of 77% using only 5% of the truncated Cora dataset for training.

3.1.1. DEPTH

Three graphs were used. Graphs 0, 200 and 400 were used. That is, the coarse-grained representation, an intermediate scale and the finest scale respectively. The number of FC hidden nodes was 30. The layer depth was varied from one layer to five layers.

The results are shown in Figure 2(a). We note that the accuracy increases initially to its highest at layer depth 2, then drops steadily as more layers are added. This is consistent with a well-known phenomenon where below a problem specific threshold, the performance of the neural network is poor, but beyond that, more depth again leads to poor performance (Loukas, 2020). It is speculated that a neural network that is too shallow is unable to learn more complex and abstract non-linear features, whereas a network that is too deep is liable to overfitting. Overly deep networks can also incur the vanishing gradient phenomenon (Ghosh & Ghosh, 2019).

3.1.2. SCALES

Each graph from the multiscale decomposition can be interpreted as a different scale of the original network. GCN depth was kept at two. The number of FC hidden nodes was 7. We note this is less than for the layer depth experiment and we note the drop in accuracy. The number of scales used, L , was varied from 3 to 176.

Figure 2(b) displays the results. There seems to be no strong trend. The highest accuracy was achieved with 8 scales. This was higher than the two lower number of scales, supporting the claim that multiscale decomposition does improve performance, but only up to a certain threshold number of scales. This is similar to the behaviour when altering the layer depth.

3.1.3. HIDDEN NODES

The GCN layer depth was kept at depth two. Three graphs were used: number 0, 200 and 400. The number of hidden nodes in the FC layer was altered from 5 to 45. The FC network always had 7 input nodes and 7 output nodes.

The results are shown in Figure 2(d). The accuracy increases sharply when increasing from 5 to 20 hidden nodes, but then a plateau is reached, where increasing the number of hidden nodes does not significantly increase the performance. Individual more extreme tests, not plotted here, were performed at 100 and 1000 hidden nodes, both still returning approximately 76% accuracy, confirming the plateau. Loukas (2020) suggests that the product of GNN width (number of hidden nodes) and depth has to exceed (a function of) the graph size to be performant. The neural network must be either deep or wide. This could explain the initial increase in accuracy as the number of hidden nodes is increased.

3.2. Noise

This experiment was conducted using standard hyperparameters: a GCN depth of two; the same three graphs from the depth and hidden node experiments; and 30 hidden nodes in the FC layer. The other experimental parameters detailed at the beginning of the hyperparameter section also apply here.

The method for noise addition selected $140p$ nodes, out of 140 in the training set, where the noise parameter, $p \in [0.1, 0.9]$. It then replaced the feature vectors corresponding to those nodes with a random vector of binary values, maintaining the original dimensions. The values were sampled from a uniform probability distribution.

The effect of varying the noise parameter on accuracy is shown in Figure 2(c). As should be expected, the accuracy decreases as the noise is increased. The model shows robustness.

4. Discussion

4.1. Limitations

There are three main limitations to the approach taken in the architecture described. First, there is added computational complexity due to the addition of a hierarchical clustering algorithm. The Girvan-Newman (2002) algorithm scales as $\mathcal{O}(N^3)$ for a sparse network. The algorithm needs to be invoked anytime a new network is input into the model. It is therefore important to choose the optimal hierarchical clustering algorithm. Second, more RAM is needed to store the graphs produced by the multiscale decomposition. A memory efficient way of storing these could be devised, such as deserializing individual disk stored graphs, from the multiscale decomposition, each time one of the scales of

GCN layers is trained, then removing them before training the next scale. Third, it is computationally more expensive to train many graph convolutional networks than training only one. The performance benefit should outweigh the extra computation time.

4.2. Future Work

Performance could be compared across multiple datasets: a model artificially generated hierarchical network, a perfectly non-hierarchical network such as a random Erdős–Rényi (1960) graph and a stochastic block model network (Karler & Newman, 2011) to understand the impact of adding communities. A neural network approach to hierarchical clustering could be implemented (Tian et al., 2014; Yang et al., 2017). Assessing performance while varying the distribution of scales, other than the linear spacing used here, may be informative.

Architecture changes could include condensing nodes in each community to one node, rather than using graphs with different numbers of edges. The decomposition would produce graphs with different numbers of nodes; each cluster represented by a node with a feature vector obtained from averaging all the node features in that cluster. They would be pushed through graph convolutional layers, then the nodes would be unpacked by expanding each node into its original cluster nodes; this time the child nodes would all have the same feature vector from the parent node. This would be performed on all graphs and would allow the graphs to have the same number of nodes again which allows the usual flattening and FC layer.

4.3. Applications

First, the architecture would be well suited for quantum molecular property prediction using the QM9 or QM7b datasets (Blum & Raymond, 2009; Montavon et al., 2013; Ramakrishnan et al., 2014; Ruddigkeit et al., 2012). This is because of its hierarchical nature as demonstrated by the improved performance of a multilevel architecture by Lu et al. (2019). Quantum molecular property prediction can aid drug discovery. Second, the protein interface prediction network by Fout et al. (2017) could be augmented with our architecture, perhaps realising performance gains. This is a prime target for our architecture due to the hierarchical nature of protein interactions.

Lastly, the graph element network by Alet et al. (2019) could be extended with our architecture. Graph element networks aim to be the neural network version of finite element analysis by using graph neural networks as computational substrates. Multiscale finite element analysis already exists where different scale meshes are used, from coarse to fine. Multiscale graph element networks do not exist, despite being a clear extension.

References

- Abu-El-Haija, S., Kapoor, A., Perozzi, B., and Lee, J. N. GCN: multi-scale graph convolution for semi-supervised node classification. *CoRR*, abs/1802.08888, 2018. URL <http://arxiv.org/abs/1802.08888>.
- Alet, F., Jeewajee, A. K., Bauza, M., Rodriguez, A., Lozano-Perez, T., and Kaelbling, L. P. Graph element networks: adaptive, structured computation and memory. *arXiv preprint arXiv:1904.09019*, 2019.
- Barabási, A.-L. and Pósfai, M. *Network science*. Cambridge University Press, Cambridge, 2016. ISBN 9781107076266 1107076269. URL <http://barabasi.com/networksciencebook/>.
- Blum, L. C. and Reymond, J.-L. 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.*, 131:8732, 2009.
- Erdős, P. and Rényi, A. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5(1):17–60, 1960.
- Fan, Y., Lin, L., Ying, L., and Zepeda-Núñez, L. A multi-scale neural network based on hierarchical matrices. *Multiscale Modeling & Simulation*, 17(4):1189–1213, 2019.
- Fout, A., Byrd, J., Shariat, B., and Ben-Hur, A. Protein interface prediction using graph convolutional networks. In *Advances in neural information processing systems*, pp. 6530–6539, 2017.
- Ghosh, S. and Ghosh, S. Exploring the ideal depth of neural network when predicting question deletion on community question answering. In *Proceedings of the 11th Forum for Information Retrieval Evaluation*, pp. 52–55, 2019.
- Girvan, M. and Newman, M. E. Community structure in social and biological networks. *Proceedings of the national academy of sciences*, 99(12):7821–7826, 2002.
- Karrer, B. and Newman, M. E. Stochastic blockmodels and community structure in networks. *Physical review E*, 83(1):016107, 2011.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. *CoRR*, abs/1609.02907, 2016. URL <http://arxiv.org/abs/1609.02907>.
- Loukas, A. What graph neural networks cannot learn: depth vs width. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=B1l2bp4YwS>.
- Lu, C., Liu, Q., Wang, C., Huang, Z., Lin, P., and He, L. Molecular property prediction: A multilevel quantum interactions modeling perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 1052–1060, 2019.
- Luzhnica, E., Day, B., and Liò, P. Clique pooling for graph classification. *CoRR*, abs/1904.00374, 2019. URL <http://arxiv.org/abs/1904.00374>.
- Montavon, G., Rupp, M., Gobre, V., Vazquez-Mayagoitia, A., Hansen, K., Tkatchenko, A., Müller, K.-R., and von Lilienfeld, O. A. Machine learning of molecular electronic properties in chemical compound space. *New Journal of Physics*, 15(9):095003, 2013. URL <http://stacks.iop.org/1367-2630/15/i=9/a=095003>.
- Ramakrishnan, R., Dral, P. O., Rupp, M., and von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1, 2014.
- Ravasz, E. and Barabási, A.-L. Hierarchical organization in complex networks. *Phys. Rev. E*, 67:026112, Feb 2003. doi: 10.1103/PhysRevE.67.026112. URL <https://link.aps.org/doi/10.1103/PhysRevE.67.026112>.
- Ruddigkeit, L., Van Deursen, R., Blum, L. C., and Reymond, J.-L. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875, 2012.
- San Kim, T., Lee, W. K., and Sohn, S. Y. Graph convolutional network approach applied to predict hourly bike-sharing demands considering spatial, temporal, and global effects. *PloS one*, 14(9), 2019.
- Tian, F., Gao, B., Cui, Q., Chen, E., and Liu, T.-Y. Learning deep representations for graph clustering. In *Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014.
- Xibin, D., Zhiwen, Y., Wenming, C., Yifan, S., and Qianli, M. A survey on ensemble learning. *Frontiers of Computer Science*, 14(2):241–258, 2020.
- Yang, C., Liu, M., Wang, Z., Liu, L., and Han, J. Graph clustering with dynamic embedding. *arXiv preprint arXiv:1712.08249*, 2017.
- Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., and Sun, M. Graph neural networks: A review of methods and applications. *CoRR*, abs/1812.08434, 2018. URL <http://arxiv.org/abs/1812.08434>.
- Zitnik, M. and Leskovec, J. Predicting multicellular function through multi-layer tissue networks. *Bioinformatics*, 33(14):i190–i198, 2017.