

BEYOND GNNs: A SAMPLE-EFFICIENT ARCHITECTURE FOR GRAPH PROBLEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite their popularity in learning problems over graph structured data, existing *Graph Neural Networks* (GNNs) have inherent limitations for fundamental graph problems such as shortest paths, k -connectivity, minimum spanning tree and minimum cuts. In all these instances, it is known that one needs GNNs of high depth, scaling at a polynomial rate with the number of nodes n , to provably encode the solution space. This in turn affects their statistical efficiency thus requiring a significant amount of training data in order to obtain networks with good performance. In this work we propose a new hybrid architecture to overcome this limitation. Our proposed architecture that we call as GNN^+ networks involve a combination of multiple parallel low depth GNNs along with simple pooling layers involving low depth fully connected networks. We provably demonstrate that for many graph problems, the solution space can be encoded by GNN^+ networks using depth that scales only *poly-logarithmically* in the number of nodes. This significantly improves the amount of training data needed that we establish via improved generalization bounds. Finally, we empirically demonstrate the effectiveness of our proposed architecture for a variety of graph problems.

1 INTRODUCTION

In recent years Graph Neural Networks (GNNs) have become the predominant paradigm for learning problems over graph structured data (Hamilton et al., 2017; Kipf & Welling, 2016; Veličković et al., 2017). Computation in GNNs is performed by each node sending and receiving messages along the edges of the graph, and aggregating messages from its neighbors to update its own embedding vector. After a few rounds of message passing, the computed node embeddings are aggregated to compute the final output (Gilmer et al., 2017). The analogy to message passing leads to a simple and elegant architecture for learning functions on graphs. On the other hand, from a theoretical and practical perspective, we also need these architectures to be *sample efficient*, i.e., learnable from a small number of training examples, where each training example corresponds to a graph. Recent works have shown that generalization in GNNs depends upon the depth of the architecture, i.e., the number of rounds of message passing, as well as the embedding size for each node in the graph (Garg et al., 2020). However, this requirement is in fundamental conflict with the message passing framework. In particular, using GNNs to compute several fundamental graph problems such as *shortest paths*, *minimum spanning tree*, *min cut* etc., necessarily requires the product of the depth of the GNN and the embedding size to scale as \sqrt{n} where n is the size of the graph (Loukas, 2020). This in turn places a significant statistical burden when learning these fundamental problems on large scale graphs. The above raises the the following question: *Can one develop sample efficient architectures for graph problems while retaining the simplicity of the message passing framework?*

Several recent works have tried to address the above question by proposing extensions to the basic GNN framework by augmenting various pooling operations in conjunction with message passing rounds to capture more global structure (Ying et al., 2018; Simonovsky & Komodakis, 2017; Fey et al., 2018). While these works demonstrate an empirical advantage over GNNs, we currently do not know of a general neural architecture that is versatile enough to *provably* encode the solution space of a variety of graph problems such as shortest paths and minimum spanning trees, while being significantly superior to GNNs in terms of statistical efficiency. In this work we propose a theoretically principled architecture, called GNN^+ networks for learning graph problems. While the basic GNN framework is inspired from classical message passing style models studied in distributed computing,

we borrow from two fundamental paradigms in graph algorithm design namely, sketching and parallel computation, to design GNN^+ networks. As a result of combining these two powerful paradigms, we get a new neural architecture that simultaneously achieve low depth and low embedding size for many fundamental graph problems. As a result our proposed GNN^+ architecture have a significantly smaller number of parameters that provably leads to better statistical efficiency than GNNs. Before we present our improved architecture, we briefly describe the standard GNN framework.

Model for GNNs. In this work we will study GNNs that fall within the message passing framework and using notation from previous works we denote such networks as GNN^{mp} (Loukas, 2020). A GNN^{mp} network operates in the AGGREGATE and COMBINE model (Gilmer et al., 2017) that captures many popular variants such as GraphSAGE, Graph Convolutional Networks (GCNs) and GIN networks (Hamilton et al., 2017; Kipf & Welling, 2016; Xu et al., 2019a). Given a graph $G = (V, E)$, let $x_i^{(k)}$ denote the feature representation of node i at layer k . Then we have

$$a_i^{(k-1)} = \text{AGGREGATE}(\{x_j^{(k-1)} : j \in N(i)\}) \quad (1)$$

$$x_i^{(k)} = \text{COMBINE}(x_i^{(k-1)}, a_i^{(k-1)}). \quad (2)$$

Here $N(i)$ is the set of neighbors for node i . Typically the aggregation and combination is performed via simple one or two layer full connected networks (FNNs), also known as multi layer perceptrons (MLPs). In the rest of the paper we will use the two terms interchangeably.

GNN^+ Networks. Our proposed GNN^+ networks consist of one or more layers of a GNN^+ block shown in Figure 1. The GNN^+ block comprises of r parallel GNN^{mp} networks followed by s parallel fully connected network modules for pooling where r and s are the hyperparameters of the architecture. More importantly we restrict the r GNN^{mp} modules to share the same set of weights. Hence the parallel GNN^{mp} modules only differ in the way the node embeddings are initialized. Furthermore, we restrict each GNN^{mp} to be of low depth. In particular, for degree- d graphs of diameter D , over n nodes, we will restrict the GNN^{mp} to be of depth $O((d + D) \cdot \text{polylog}(n))$. Similarly, we require the s fully connected networks to be of depth $O((d + D) \cdot \text{polylog}(n))$ and share the network weights. We connect the outputs of the GNN^{mp} modules to the fully connected pooling networks in a sparse manner and restrict the input size of each fully connected network to be $O((d + D) \cdot \text{polylog}(n))$. Stacking up L layers of GNN^+ blocks results in a GNN^+ network that is highly parameter efficient and in total has $O((d + D)L \cdot \text{polylog}(n))$ parameters. For such a network we call the depth as the total number of message passing rounds and the number of MLP layers used across all the L stacks. Since we restrict our MLPs and GNN^{mp} blocks inside a GNN^+ network to be of low depth, we will often abuse notation and refer to a GNN^+ architecture with L stacks of GNN^+ blocks as a depth L architecture. Our proposed design lets us alternate between local computations involving multiple parallel GNN blocks and global post-processing stages, while still being sample efficient due to the enforced parameter sharing. We will show via several applications that optimal or near-optimal solutions to many popular graph problems can indeed be computed via a GNN^+ architecture. Below we briefly summarize our main results.

Overview of Results. To demonstrate the generality of our proposed GNN^+ architecture, we study several fundamental graph problems and show how to construct efficient GNN^+ networks to compute optimal or near optimal solutions to these problems. In particular, we will focus on degree- d graphs, i.e., graphs of maximum node degree d , with n nodes and diameter D and will construct GNN^+ networks of depth $\text{polylog}(n)$ and $O((D + d)\text{polylog}(n))$ total parameters.

Shortest Paths. The first problem we consider is the fundamental graph problem of computing (approximate) all pairs shortest paths in undirected graphs. Given a graph $G = (V, E)$, let $d_G(u, v)$ be the shortest path between nodes u and v . We say that an output $\{\hat{d}_G(u, v) : u, v \in V\}$ is an

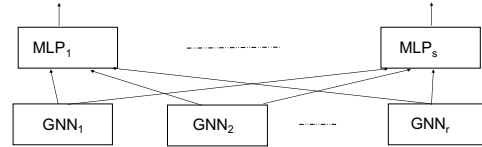


Figure 1: The basic GNN^+ block.

α -approximate solution if for all $u \neq v$ it holds that

$$d_G(u, v) \leq \tilde{d}_G(u, v) \leq \alpha d_G(u, v).$$

We construct efficient GNN^+ networks for all pairs shortest paths with the following guarantee.

Theorem 1 (Informal Theorem). *For any constant $c > 1$, there is a depth $O(D \log d + \log n)$ GNN^+ network with $O((n^{\frac{2}{c}} + d) \text{polylog}(n))$ parameters that computes $(4c - 2)$ -approximate all pairs shortest paths in the undirected unweighted degree- d graphs over n nodes. On the other hand, computing a c -approximate shortest paths using GNN^{mp} networks requires a network of depth $\Omega(n)$.*

From the above theorem we see that by setting $c = O(\log n)$ we can encode a c -approximate solution using an $O(D \log d + \log n)$ GNN^+ network with only $O(d \cdot \text{polylog}(n))$ parameters. This is in stark contrast with the depth requirement of $\Omega(n)$ for the traditional GNN^{mp} networks.

Connectivity Measures. Next we consider computing various graph connectivity measures. We first study the popular measure based on graph effective resistances (Chandra et al., 1996).

Definition 1 (Effective Resistance). *Let G be a weighted undirected graph G with adjacency matrix A and the associated Laplacian $L = D - A$. Given an edge u, v , the effective resistance between u, v is defined as*

$$R_{u,v} = \xi_{u,v}^\top L^\dagger \xi_{u,v}.$$

Here $\xi_{u,v}$ is an n dimensional vector with $+1$ at position u , -1 at position v and zeros everywhere. L^\dagger refers to the matrix pseudo-inverse.

We also study the following connectivity measure that was proposed by Panigrahy et al. (2012) in the context of web graphs. Given an undirected graph G , let G_p be the random graph obtained by sampling each edge with probability p .

Definition 2 (Affinity). *For any two vertices u, v and for $p \in [0, 1]$, define $A_p(u, v)$ to be the probability that u, v are connected in G_p . Then the affinity between u and v is defined as*

$$A(u, v) = \mathbb{E}_p[A_p(u, v)]$$

where the expectation is taken over p drawn from the uniform distribution in $[0, 1]$.

For the above measures we show the following

Theorem 2 (Informal Theorem). *There exists a GNN^+ architecture with $O(D \log(nd))$ parameters, and depth $O(D \log(nd))$ on graphs of diameter D with n nodes and maximum degree d , that approximate the above connectivity measures up to constant factors. On the other hand using GNN^{mp} networks to compute the above measures, even approximately, necessarily requires a network of depth $\Omega(\sqrt{n})$.*

Clustering, Minimum Cuts and Minimum Spanning Trees. Finally, we showcase the power of a GNN^+ architecture for computing other fundamental graph problems. Given an undirected graph G , the spectral clustering of G corresponds to the cut obtained by taking the sign of the eigenvector v corresponding to the second smallest eigenvalue $\lambda_2(L)$, where L is the graph Laplacian. For computing the spectral clustering via GNN^+ networks we show the following

Theorem 3 (Informal Theorem). *There is a GNN^+ network of depth $\ell = O(\frac{1}{\lambda_2(L)\epsilon^2} \log n)$, with $O(d)$ parameters that computes an ϵ -approximate spectral clustering on graphs of degree d . On the other hand, using GNN^{mp} networks to even approximately compute the spectral clustering requires depth $\Omega(\sqrt{n})$.*

Next we consider the classical problems of computing a global minimum cut and minimum spanning trees in undirected graphs.

Theorem 4 (Informal Theorem). *There exist GNN^+ networks of of depth $O((D + \log n) \log n)$, and $O(d)$ parameters for computing a global minimum cut (MINCUT) and minimum spanning tree (MST) in degree d graphs of diameter D . Furthermore, using GNN^{mp} networks to compute these primitives (even approximately) necessarily requires depth $\Omega(\sqrt{n})$.*

Generalization Bounds. Our final result concerns the generalization properties of a depth L GNN^+ architecture. For ease of exposition, we state here the results for the case when the GNN^+ architecture produces a one dimensional output. More general results are presented in Appendix D. Our generalization bounds depend on the depth L and the total number of parameters P in the GNN^+ network. Following recent work on providing generalization bounds for fully connected and convolutional neural networks (Bartlett et al., 2017; Long & Sedghi, 2020) that are based on *distance to initialization*, we consider the class \mathcal{F}_β of depth L GNN^+ networks with P parameters that are at a distance β from a reference parameter configuration (typically the parameters at random initialization). Let $y \in \mathbb{R}$ denote the output of the network and consider a Lipschitz loss function $\ell(y, \hat{y})$. Then, we provide following guarantee.

Theorem 5 (Informal Theorem). *Let $\ell(\hat{y}, y)$ be a Lipschitz loss function bounded in $[0, B]$. Then, given m i.i.d. samples $(G_1, y_1), (G_2, y_2), \dots (G_m, y_m)$ generated from a distribution D , with probability at least $2/3$, it holds that for all $f \in \mathcal{F}_\beta$,*

$$\left| \hat{\mathbb{E}}_D[\ell_f] - \mathbb{E}_D[\ell_f] \right| \leq O\left(B\sqrt{\frac{P(\beta + L)}{m}}\right).$$

We refer the reader to Theorem 16 in Appendix D for a formal statement and the proof. Notice that the above theorem implies that our proposed GNN^+ architecture for the above graph problems can indeed be trained using very few samples as opposed to the traditional GNN^{mp} networks since the GNN^+ network requires much fewer parameters and depth. Furthermore, since a GNN^{mp} network is a special case of a GNN^+ architecture, our analysis also leads to an improved bound on the generalization guarantees for GNN^{mp} networks as well. In particular, the above theorem improves upon the recent work of Garg et al. (2020) that provides generalization guarantees for training GNNs that scale with the branching factor of the graph. Using our improved analysis we are able to remove this dependence on the branching factor. See Appendix D for details.

2 RELATED WORK

GNNs operate primarily in the message passing framework where nodes aggregate and combine messages from their neighbors to update their embeddings. Several variants of this basic paradigm have been proposed, with each differing in how the aggregation and combination is performed. Popular variants include GraphSAGE (Hamilton et al., 2017), Graph Convolutions Networks (Kipf & Welling, 2016), GIN networks (Xu et al., 2019a), and graph pooling (Ying et al., 2018).

Various recent works have also studied the representation power of GNNs. The work of Xu et al. (2019a) demonstrates that the GNNs as considered in equation 1 are as powerful as the Weisfeiler-Lehman test for graph isomorphism (Weisfeiler & Lehman, 1968). The recent work of Xu et al. (2019b) compares the message passing framework of GNNs in representing computations involving dynamic programming. GNN networks that can capture higher order variants of the WL test have also been proposed recently (Maron et al., 2019).

Several works have also explored the limitations of GNNs for computing graph primitives. The work of Loukas (2020) established a correspondence between the message passing GNN framework and the well studied CONGEST model of distributed computing (Peleg, 2000). Based on the above correspondence it follows that in order to represent several important graph problems such as shortest paths, minimum cuts and minimum spanning tree, either the depth of the GNN or the embedding size of the nodes has to scale with the graph size at a polynomial rate. Notice that these lower bounds apply to any form of message passing framework and as a result recent work in incorporating non-symmetric node messages (Sato et al., 2019) in GNNs also run into the same barriers.

In order to address the above limitations recent works have proposed combining the GNN architecture with pooling mechanisms for aggregating more global information (Ying et al., 2018; Defferrard et al., 2016; Simonovsky & Komodakis, 2017; Fey et al., 2018; Bianchi et al., 2019; Du et al., 2019). For example the work of Ying et al. (2018) proposes a hierarchical approach where a GNN network is followed by a clustering step to compute higher level “nodes” to be used in the subsequent GNN operation. While these approaches show empirical promise, ours is the first work to design a principled architecture with theoretical guarantees that merges local distributed computations with global postprocessing stages.

Finally, the question of generalization for GNNs has also been studied in recent works. The most relevant to us is the recent work of Garg et al. (2020) that analyzes the Rademacher complexity of GNNs with the aggregate mechanism being a simple addition and the combine mechanism being a one layer neural network. Via analyzing the Rademacher complexity the authors show that the generalization for GNNs depends on the depth, the embedding size and the branching factor of the graph. Our improved analysis in Section D extends the result of Garg et al. (2020). Not only does our generalization bound apply to the more general GNN^+ networks, for the case of GNNs considered in (Garg et al., 2020) our analysis shows that the dependence on the branching factor can be eliminated in the generalization bounds. Generalization bounds have also been proved recently for GNN based networks that use the Neural Tangent Kernel (NTK) during the aggregation and combination operations (Du et al., 2019).

3 SHORTEST PATHS

In this section we provide a proof sketch of Theorem 1 showing how to construct an efficient GNN^+ architecture for the Shortest Paths problem. In particular we study all pairs shortest paths.

All Pairs Shortest Paths. The input is a graph $G = (V, E)$ with n nodes. The desired output is an $\binom{n}{2}$ dimensional vector containing (approximate) shortest path values between each pair of vertices. Given a graph G , let $d_G(u, v)$ be the shortest path between nodes u and v . We say that an output $\{\tilde{d}_G(u, v) : u, v \in V\}$ is an α -approximate solution if for all $u \neq v$ it holds that

$$d_G(u, v) \leq \tilde{d}_G(u, v) \leq \alpha d_G(u, v).$$

We first show that the GNN^{mp} networks are highly inefficient for learning this problem.

Theorem 6. *Consider a GNN^{mp} \mathcal{N} of depth L over n nodes where each node has a representation size of B . If \mathcal{N} encodes α -approximate all pairs shortest paths for graphs of diameter bounded by D , and for $\alpha < 3/2$, then it must hold that $B \cdot L \geq \Omega(n)$. Furthermore, the lower bound holds for undirected unweighted graphs as well.*

Proof. The recent work of Loukas (2020) established that computation in GNN^{mp} networks is equivalent to the CONGEST model of computation popularly studied in the design of distributed algorithms (Peleg, 2000). In particular, a lower bound on the product of depth (L) and representation size (B) can be obtained by establishing the corresponding lower bound on the product of the number of rounds and the size of messages in the CONGEST model of computing. Furthermore, the result of Holzer & Wattenhofer (2012) shows that in the CONGEST model approximating all pairs shortest paths, even on unweighted undirected graphs requires the product of the number of rounds and the message size to be $\Omega(n)$. Hence the corresponding lower bound on $B \cdot L$ follows. \square

Circumventing Lower Bounds via GNN^+ . Next we detail our proposed GNN^+ architecture that can encode approximate shortest paths with significantly smaller depth and parameter requirements.

Unweighted Graphs. To illustrate the main ideas we study the case of undirected unweighted graphs. See Appendix A for the more general case of weighted graphs. The starting point of our construction is the following fundamental theorem of Bourgain (1985) regarding metric embeddings.

Theorem 7 ((Bourgain, 1985)). *Any n -point metric (X, d) can be embedded into the Euclidean metric of dimensionality $O(\log n)$ and distortion $O(\log n)$.*

The above theorem suggests that in principle, if we only want to estimate shortest paths up to an approximation of $O(\log n)$, then we only need node embeddings of size $O(\log n)$. If there were a GNN^{mp} network that could produce such embeddings, then one could simply compute the Euclidean distance between each pair of points to get the approximate shortest path. Furthermore, computing the Euclidean distance given the node embeddings can be done easily via a low depth full connected network. Unfortunately, producing the necessary low dimensional embeddings is exactly the task for which GNN^{mp} networks require large depth as proved in Theorem 6 above. While there do exist semi-definite programming based algorithms (Linial et al., 1995) for computing the embeddings required for Bourgain’s theorem, they are not suitable for implementation via efficient neural architectures. Instead we rely on sketching based algorithms for computing shortest path distances.

In particular, for the unweighted case we adapt the sketch based approximate shortest path algorithms of Das Sarma et al. (2010) for designing an efficient network architecture. The sketch proposed in the work of Das Sarma et al. (2010) computes, for each node u , the distance of u from a random subset S of the nodes. This can be done via a simple breadth first search (BFS). Repeating this process k -times provides a k -dimensional embedding for each vertex and for an appropriate choice of k , these embeddings can be used to compute approximate shortest paths. Notice that this sketch based procedure is highly amenable to be implemented in a message passing framework. Overall, the algorithm performs multiple parallel BFS subroutines to compute the embeddings. It is also well known that BFS on diameter D can be implemented by a GNN^{mp} of depth $O(D)$.

Based on the above intuition, our proposed architecture is shown in Figure 2. It consists of k parallel breadth first search (BFS) modules for $k = \Theta(n^{\frac{1}{c}} \log n)$ for a constant $c > 1$. Module i computes the shortest path from each vertex in G to any vertex in the set S_i . The sets S_1, S_2, \dots, S_k are randomly chosen subsets of the vertex set V of various sizes. In particular there are $\Theta(n^{\frac{1}{c}})$ subsets of size 1, $\Theta(n^{\frac{1}{c}})$ subsets of size 2, $\Theta(n^{\frac{1}{c}})$ subsets of size 2^2 , and so on up to $\Theta(n^{\frac{1}{c}})$ subsets of size $2^{\lceil \log n \rceil}$. The BFS module i produces n distance values $v_1^{(i)}, \dots, v_n^{(i)}$. These modules are followed by $\binom{n}{2}$ fully connected networks where each module is responsible for computing the approximate shortest path distance between a pair of vertices. In particular we have $\tilde{d}_G(s, t) = \max_i |v_s^{(i)} - v_t^{(i)}|$.

Notice from the discussion in Section 1 that the architecture in Figure 2 is a GNN^+ network with a single GNN^+ block. In the next section we will show how we can generalize to a suite of graph problems by stacking up multiple GNN^+ blocks. For our proposed network we have the following guarantee.

Theorem 8. *For any constant $c > 1$, and for a fixed graph topology over n nodes with maximum degree d and diameter D , there exists a neural network as shown in Figure 2 of size $O(n^{2+1/c})$, $\tilde{O}(n^{\frac{2}{c}})$ parameters, and depth $O(D \log d + \log n)$, that encodes $(2c - 1)$ -approximate all pairs shortest paths in G .*

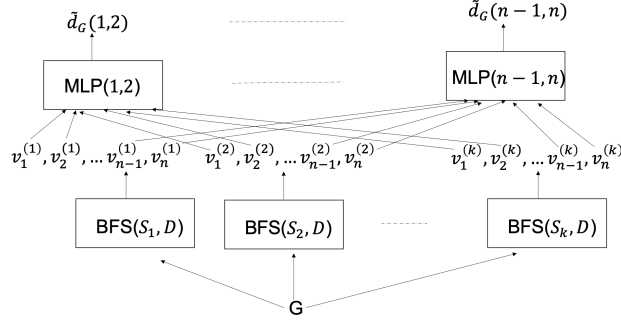


Figure 2: The network architecture for approximate all pairs shortest paths in unweighted graphs.

Before proving the theorem above we establish two supporting lemmas concerning the implementation of the BFS modules and the MLP module in the network architecture described in Figure 2.

Lemma 1. *The BFS module in Figure 2 can be implemented by a GNN of depth $O(D)$, $O(1)$ total parameters and with each node having a representation size of $O(1)$.*

Lemma 2. *For any k , the MLP module in Figure 2 can be implemented by a network of depth $O(\log k)$, $O(k^2)$ total parameters.*

Proof of Theorem 8. The correctness of the network architecture follows from the work of Das Sarma et al. (2010). Next we establish bounds on the total depth, size and the number of parameters. We have $k = \Theta(n^{\frac{1}{c}} \log n)$ copies of the BFS module. Each BFS module is of size $O(nd \log d)$ since there are n nodes and each node implements a min function of size $O(d \log d)$. Hence, in total the BFS modules have size $O(n^{1+1/c} d \log d \log n)$. Next we have $\binom{n}{2}$ MLP modules each of size $O(k \log k)$ for a total size of $O(n^{2+1/c} \log n)$. Hence the total size of the neural network is bounded by $O(n^{2+1/c} \log n)$.

Next we bound the depth and the total number of parameters. The BFS module has $O(D)$ rounds with each requiring a depth of $O(\log d)$ for a total depth of $O(D \log d)$. The MLP module has a depth bounded by $O(\log k) = O(\log n)$. Hence the total depth is $O(D \log D + \log n)$. Finally, the BFS module requires $O(1)$ parameters and the MLP module requires $O(k^2)$ parameters. Hence, the total number of parameters in our architecture are bounded by $O(k^2) = O(n^{2/c})$. \square

4 MINIMUM CUTS

To illustrate another application, in this section we design an efficient GNN^+ based architecture for computing the minimum cut in an undirected graph. We first argue in Appendix C that even computing an approximate mincut using traditional GNN^{mp} networks requires $\Omega(\sqrt{n})$ rounds. Our efficient GNN^+ based architecture is based on the parallel algorithm for computing mincut (Karger & Stein, 1996) and is shown in Figure 3. More importantly the architecture comprises of multiple layers of GNN^+ blocks in contrast to a single GNN^+ block in the case of shortest paths. The algorithm of Karger & Stein (1996) relies on the following lemma.

Lemma 3 ((Karger & Stein, 1996)). *Let m be the number of edges in the graph and let L be a random ordering of the m edges. Then with probability at least $\frac{1}{n^2}$, there exists a prefix L' of L such that contracting the edges in L' results in exactly two connected components corresponding to the global minimum cut. in the graph.*

The above suggests a natural neural architecture of depth $O(\log m) = O(\log n)$ that aims to mimic a binary search procedure to infer the correct prefix L' . At each of the $O(\log n)$ stages, each vertex maintains a list of which of its connecting edges are *active*, i.e., not contracted. This requires $O(d)$ representation size. The goal next is to infer whether the number of connected components induced by the active edges is one, two, or more than two. This in turn decides the part of the ordering the next stage will focus on and which edges will become inactive during the current stage. The computation of connected components can be carried out using at most two breadth first searches and hence via $O(D)$ rounds of a GNN^{mp} network. Following this intuition we arrive at the proposed architecture in Figure 3. Formally, we have the following guarantee.

Theorem 9. *For a fixed graph topology over n nodes with maximum degree d and diameter D , the network in Figure 3 produces the minimum cut. Furthermore, the network is of depth $\ell = O(D \log^2 n)$, size $O(n\ell)$, and has $O(d + \log n)$ parameters.*

Proof. Each vertex maintains an $O(d)$ sized representation indicating which of its edges are currently *active*. The GNN module simply performs a BFS over active edges starting from a specified vertex (output by the *Update Prefix* module). Initially random ordering of the edges is chosen. The goal of the *Update Prefix* module is to update the current prefix in this ordering that specifies the set of active edges. Let L be the initial ordering of the edges. Then the guarantee from Karger & Stein (1996) implies that with probability at least $1/n^2$ there exists a prefix L' of L such that the subgraph comprising of edges in L' consists of exactly two connected components corresponding to the optimal mincut. Hence, given the output of the BFS performed by the GNN, the *Update Prefix* module has to decide if the number of connected components is more than two, and update the prefix accordingly. This can be implemented by an MLP using $O(\log n)$ depth and $O(\log m)$ parameters where m is the number of edges in the graph. Hence overall we get $O(D \log n)$ depth and $O(d + \log n)$ parameters. \square

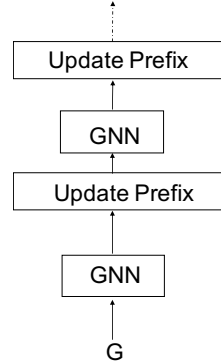


Figure 3: The network architecture for minimum cut.

5 EXPERIMENTS

We show the efficacy of GNN^+ on the aforementioned graph problems: Shortest Paths, Effective Resistance, Affinity, MINCUT and MST, and compare to a state-of-the-art GNN^{mp} model (Xu et al., 2019a).

Dataset. We generated synthetic random graphs between 500 and 1000 nodes. For the affinity measure, we used graphs with 250 nodes because of the need for using very dense graphs to have a reasonable number of alternate paths between any two end points. In general, we generated the data sets as follows: we fix the number of nodes n in the graph to take values in $[250, 1000]$. For each value of n we generate graphs from the Erdos-Renyi model $G(n, p)$ with edge sampling probability $p = \frac{\alpha}{n}$.

We vary α depending on the problem. Specifically, we set α to be a constant in $[1, 100]$ to capture varying degrees of sparsity. For each n, p we generate 30,000 training examples consisting of tuples of the form $(G, s, t, d(s, t))$ where G is a random graph drawn from $G(n, p)$, s, t are two vertices uniformly drawn at random and $d(s, t)$ is one of shortest path value, effective resistance, or affinity between the two vertices. In the case of min cut and minimum spanning tree, we generate tuples (g, v_G) where v_G corresponds to the value of the minimum spanning tree or the global minimum cut.

Models and Configurations. For our baseline GNN^{mp} implementation, we used the GIN model proposed in Xu et al. (2019a). This has been empirically shown (Xu et al., 2019a; Loukas, 2020; Errica et al., 2020) to be a state-of-the-art GNN^{mp} model on several datasets. GIN updates feature representations $x_v^{(k)}$ of each node v at iteration k as: $x_v^{(k)} = \text{MLP}\left((1 + \epsilon^{(k)}) \cdot x_v^{(k-1)} + \sum_{u \in N(v)} x_u^{(k-1)}\right)$, where MLP refers to a Multi-Layer Perceptron, $N(v)$ is the set of neighbors of v , and ϵ is a learnable parameter. For problems that involved weighted graphs (e.g. MST), we incorporated edge weights into the GIN update equation by replacing the sum of neighbor representations by a weighted sum.

Our GNN^+ implementation also used the same GIN implementation as its internal GNN^{mp} block. All graphs in our experiments were undirected. For both baseline and GNN^+ , we used node degree as the input node features for MINCUT and MST. For Shortest Paths, Effective Resistance and Affinity, we set input node features to be Booleans indicating if the node is a source/destination node or not.

Following Xu et al. (2019a), we performed 10-fold cross-validation for each of our experiments (corresponding to the two models and five problems), and report the average validation mean squared error (MSE) across the 10 folds. We run each 10-fold cross-validation experiment 10 times to compute confidence intervals. We apply batch normalization at each layer, and use an Adam optimizer and decay the learning rate by 0.5 every 50 epochs, and train for up to 600 epochs. For both the baseline GNN^{mp} and the GNN^+ model, we tune the following parameters: initial learning rate $\in \{0.001, 0.003, 0.005, 0.007, 0.01, 0.03, 0.05\}$, number of hidden units $\in \{8, 16, 32, 64\}$, batch-size $\in \{32, 64\}$, dropout $\in \{0, 0.5\}$. For GNN^{mp} we also tuned the depth (number of layers) $\in \{2, 4, 8, 12\}$. For the GNN^+ model, we tuned the number of parallel GNNs in each GNN^+ block to $\in \{1, 2, 3\}$ with GNN depth $\in \{2, 4\}$. We also tuned the number of GNN^+ layers $\in \{1, 2, 3\}$. We fixed the depth of each MLP block in GNN^{mp} and GNN^+ to 2.

Problem	Label Variance	Avg. MSE (GNN^{mp})	Avg. MSE (GNN^+)
Shortest Path	7.482	0.975 ± 0.031	0.849 ± 0.022
Effective Resistance	7.949	0.397 ± 0.025	0.187 ± 0.008
Affinity	3.030	$0.0025 \pm 1.75\text{e-}04$	$0.0018 \pm 1.89\text{e-}05$
MST	4637.4	1011.39 ± 106.94	733.901 ± 30.97
MINCUT	11.964	0.963 ± 0.110	0.694 ± 0.07

Table 1: Performance of the GNN^{mp} and GNN^+ architectures

Results. To validate our theory regarding the better generalization bounds for GNN^+ models compared to GNN^{mp} , we compare the test mean squared errors for the two models across the five datasets. For all the five problem, Table 1 lists the test MSEs and corresponding standard deviations for the two models. As a sanity check, we also plot the variance of the labels in our datasets, which corresponds to the MSE obtained by a naive model that predicts the mean label. We observe significant gains in accuracy of anywhere between 15% relative MSE improvement over the GNN^{mp} baseline (for Shortest Paths) to as much as 108% relative MSE improvement (for Effective Resistance). Note that the naive mean predictor’s MSE is at least an order of magnitude larger than all the MSE values for GNN^{mp} and GNN^+ (except for the MST dataset, where it is around five times larger - we suspect that the weighted graphs involved in this dataset make this a harder problem).

We posit that these accuracy gains directly stem from the sample-efficiency of the GNN^+ models as captured in Theorems 1, 2 and 4 - the most compact GNN^+ networks that can represent these problems are smaller than the corresponding most compact GNN^{mp} networks. Hence, by Theorem 5, such networks will have smaller generalization errors. In the appendix, we also plot the test accuracy as a function of number of epochs that suggest that our models also converge faster than the baseline GNN^{mp} models, though we do not have any theoretical justification supporting this observation.

REFERENCES

- Raman Arora, Amitabh Basu, Poorya Mianjy, and Anirbit Mukherjee. Understanding deep neural networks with rectified linear units. *arXiv preprint arXiv:1611.01491*, 2016.
- László Babai, Peter Frankl, and Janos Simon. Complexity classes in communication complexity theory. In *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*, pp. 337–347. IEEE, 1986.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, pp. 6240–6249, 2017.
- Filippo Maria Bianchi, Daniele Grattarola, and Cesare Alippi. Mincut pooling in graph neural networks. *arXiv preprint arXiv:1907.00481*, 2019.
- Jean Bourgain. On lipschitz embedding of finite metric spaces in hilbert space. *Israel Journal of Mathematics*, 52(1-2):46–52, 1985.
- Ashok K Chandra, Prabhakar Raghavan, Walter L Ruzzo, Roman Smolensky, and Praseen Tiwari. The electrical resistance of a graph captures its commute and cover times. *Computational Complexity*, 6(4):312–340, 1996.
- Atish Das Sarma, Sreenivas Gollapudi, Marc Najork, and Rina Panigrahy. A sketch-based distance oracle for web-scale graphs. In *Proceedings of the third ACM international conference on Web search and data mining*, pp. 401–410, 2010.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems*, pp. 3844–3852, 2016.
- Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, pp. 5723–5733, 2019.
- Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph neural networks for graph classification. *ICLR*, 2020.
- Matthias Fey, Jan Eric Lenssen, Frank Weichert, and Heinrich Müller. Splinecnn: Fast geometric deep learning with continuous b-spline kernels. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 869–877, 2018.
- Sebastian Forster and Danupon Nanongkai. A faster distributed single-source shortest paths algorithm. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 686–697. IEEE, 2018.
- Vikas K. Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits of graph neural networks. *ICML*, 2020.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*, 2017.
- Evarist Giné and Armelle Guillou. On consistency of kernel density estimators for randomly censored data: rates holding uniformly over adaptive intervals. In *Annales de l’IHP Probabilités et statistiques*, volume 37, pp. 503–522, 2001.
- Will Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pp. 1024–1034, 2017.
- Stephan Holzer and Roger Wattenhofer. Optimal distributed all pairs shortest paths and applications. In *Proceedings of the 2012 ACM symposium on Principles of distributed computing*, pp. 355–364, 2012.
- David R. Karger and Clifford Stein. A new approach to the minimum cut problem. *Journal of the ACM*, 43:601–640, 1996.

- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Nathan Linial, Eran London, and Yuri Rabinovich. The geometry of graphs and some of its algorithmic applications. *Combinatorica*, 15(2):215–245, 1995.
- Philip M. Long and Hanie Sedghi. Generalization bounds for deep convolutional neural networks, 2020.
- Andreas Loukas. What graph neural networks cannot learn: depth vs width. *ICLR*, 2020.
- Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. In *Advances in Neural Information Processing Systems*, pp. 2156–2167, 2019.
- Rina Panigrahy, Marc Najork, and Yinglian Xie. How user behavior is related to social affinity. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pp. 713–722, 2012.
- David Peleg. *Distributed computing: a locality-sensitive approach*. SIAM, 2000.
- Atish Das Sarma, Stephan Holzer, Liah Kor, Amos Korman, Danupon Nanongkai, Gopal Pandurangan, David Peleg, and Roger Wattenhofer. Distributed verification and hardness of distributed approximation. *SIAM Journal on Computing*, 41(5):1235–1265, 2012.
- Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Approximation ratios of graph neural networks for combinatorial problems. In *Advances in Neural Information Processing Systems*, pp. 4081–4090, 2019.
- Martin Simonovsky and Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3693–3702, 2017.
- M Sollin. La trace de canalisation. *Programming, Games, and Transportation Networks*, 1965.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Boris Weisfeiler and Andrei A Lehman. A reduction of a graph to a canonical form and an algebra arising during this reduction. *Nauchno-Tekhnicheskaya Informatsia*, 2(9):12–16, 1968.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *ICLR*, 2019a.
- Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019b.
- Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Advances in neural information processing systems*, pp. 4800–4810, 2018.

A SHORTEST PATHS

Unweighted Graphs. We first provide more technical details for the case of unweighted graphs, followed by the generalization to the weighted case.

We start with describing the construction of the BFS module in Figure 2. Recall that the BFS module i produces n distance values $v_1^{(i)}, \dots, v_n^{(i)}$. These modules are followed by $\binom{n}{2}$ fully connected networks where each module is responsible for computing the approximate shortest path distance between a pair of vertices. In particular we have

$$\tilde{d}_G(s, t) = \max_i |v_s^{(i)} - v_t^{(i)}|.$$

We restate and prove the following lemmas concerning the implementation of the BFS modules and the MLP module in the network architecture described in Figure 2.

Algorithm 1 BFS module

```

1: function BFS(Graph  $G$ , subset  $S$ )
2:   Initialize  $x_i^{(0)} = 0$  if  $i \in S$  and  $\infty$  otherwise.
3:   for  $t = 1, \dots, D$  do
4:
```

$$a_i^{(k-1)} = \min(\{x_j^{(k-1)} : j \in N(i)\})$$

$$x_i^{(k)} = \min(x_i^{(k-1)}, a_i^{(k-1)}).$$

```

5:   Output  $\{x_1^{(D)}, \dots, x_n^{(D)}\}$ .
```

Lemma 4 (Restatement of Lemma 1). *The BFS module in Figure 2 can be implemented by a GNN of depth $O(D)$, $O(1)$ total parameters and with each node having a representation size of $O(1)$.*

Proof. The implementation of the BFS module is shown in Algorithm 1. Notice that each node i stores two values, namely a_i^ℓ and v_i^ℓ where ℓ denotes the current depth of the GNN. In total we have a depth of D corresponding to the total number of iterations and in each iteration one needs to compute the minimum of up to d values. The min operation can be implemented by a one dimensional CNN of depth $O(\log d)$ with a kernel of size 2 and ReLU activations, as described in Arora et al. (2016). Since the min operation is repeatedly applied and has a kernel of constant size, in total we have $O(1)$ parameters. \square

Next we show how to implement the MLP module in Figure 2.

Algorithm 2 MLP module

```

1: function MLP( $v_s^{(1)}, v_t^{(1)}, \dots, v_s^{(k)}, v_t^{(k)}$ )
2:   For  $i = 1, \dots, k$ , compute  $r_{i,1} = v_s^{(i)} - v_t^{(i)}$  and  $r_{i,2} = v_t^{(i)} - v_s^{(i)}$ .
3:   Output  $\max_i \max_j r_{i,j}$ .
```

Lemma 5 (Restatement of Lemma 2). *For any k , the MLP module in Figure 2 can be implemented by a network of depth $O(\log k)$, $O(k^2)$ total parameters.*

Proof. The MLP module is sketched in Algorithm 2. In the first step the module computes $r_{i,j}$ values for $i = 1$ to k and $j \in \{1, 2\}$. This can be easily computed by a depth one linear network of width $2k$ and $O(k^2)$ total parameters. In the next step, the module has to compute the maximum of $2k$ values. Similar to the implementation of the BFS module, this can be done by a depth $O(\log k)$ network with ReLU activations and $O(1)$ total parameters. \square

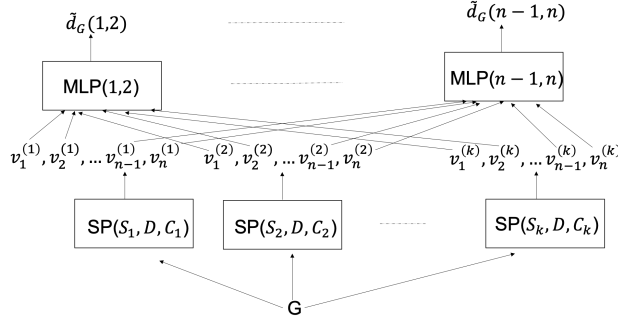


Figure 4: The network architecture for approximate all pairs shortest paths in weighted graphs.

Weighted Graphs. We next consider the case of weighted graphs with polynomially bounded weights. We denote by $w(u, v)$ the weight of the edge $(u, v) \in E$. In this case, we again follow the sketch based approach of Das Sarma et al. (2010) as in Figure 2. However, we replace the BFS module with a shortest path module to compute weighted shortest paths from vertices in sets S_i to the rest of the vertices in G . The new architecture in this case is shown in Figure 4. Next we

Algorithm 3 SP module

- 1: **function** SP(G, S, C)
- 2: Set $h = |C|$. Compute h hop distances $\hat{d}(x, v)$ from each $x \in C$ to each $v \in V$.
- 3: Using distances \hat{d} compute shortest distance $d'()$ between vertices in C to any vertex in S .
- 4: Compute and output h -hop distances between vertices in G to any vertex in S using the weighting scheme w' where

$$w'(u, v) = \begin{cases} d'(u, v), & \text{if } (u, v) \in (S \times C) \setminus E \\ 2w(u, v), & \text{if } (u, v) \in E \setminus (S \times C) \\ \min(d'(u, v), 2w(u, v)), & \text{if } (u, v) \in E \cap (S \times C) \end{cases}$$

describe the implementation of the shortest path module (SP). The module implements the algorithm of Forster & Nanongkai (2018) who designed a 2-approximate single source shortest path algorithm in the CONGEST model of computation. As before SP module i takes as input the graph G , the set of vertices S_i . Furthermore, it also takes as input a subset C_i of $O(\sqrt{nD})$ vertices that are used to approximate shortest path distances. The SP module is sketched in Algorithm 3. We next prove the following regarding implementing the SP module as a neural network.

Lemma 6. *The Shortest Path module in Algorithm 3 can be implemented by a neural network with depth $O(\sqrt{nD} \log d)$, size $O(n\sqrt{nD})$ and $O(1)$ total parameters.*

Proof. We will implement the shortest path module using the GNN framework. Each node will have a representation size of $O(d + \log n)$. This will let the hold original graph distances $w()$, as well as the distances \hat{d} and the weights w' in step 4 of the algorithm. Step 2 of the algorithm can be implemented by using $|C|$ parallel GNN modules with each implementing the Bellman-Ford algorithm up to $h = \sqrt{nD}$ rounds. As in the BFS implementation, this requires $O(1)$ parameters and $O(\sqrt{nD} \log d)$ total depth. Step 3 can be similarly implemented by running Bellman-Ford algorithm for $|C| = \sqrt{nD}$ steps. This adds a depth of $O(\sqrt{nD} \log d)$ to the network. Finally, step 4 can again be performed by $O(\sqrt{nD})$ rounds of Bellman-Ford. Hence, in total we have $O(\sqrt{nD} \log d)$ depth and $O(1)$ parameters. \square

Using the lemma above and the analysis of Theorem 8 we get the following result for weighted graphs.

Corollary 1. *For any constant $c > 1$, and for any graph over n nodes with maximum degree d and diameter D , there exists a neural network as shown in Figure 4 of size $O(n^{2+1/c})$, $\tilde{O}(n^{\frac{2}{c}})$ parameters,*

and depth $O(\sqrt{nD} \log d + \log n)$, that encodes $(4c - 2)$ -approximate all pairs shortest paths in the weighted graph G .

B COMPUTING CONNECTIVITY MEASURES

In this section we study the limitations of using GNNs for computing various graph connectivity measures and show how the GNN^+ architecture can overcome these. We first study the following connectivity measure as proposed in the work of Panigrahy et al. (2012). Given an undirected graph G , let G_p be the random graph obtained by sampling each edge with probability p .

Definition 3 (Affinity). *For any two vertices u, v and for $p \in [0, 1]$, define $A_p(u, v)$ to be the probability that u, v are connected in G_p . Then the connectivity measure between u and v is defined as*

$$A(u, v) = \mathbb{E}_p[A_p(u, v)]$$

where the expectation is taken over p drawn from the uniform distribution in $[0, 1]$.

The work of Panigrahy et al. (2012) showed that the above measure is closely related to the notion of strong connectivity in a graph. We will also consider the following weighted version of the above measure.

Definition 4. *Given a weighted undirected graph $G = (V, E, w)$ with edge weights in $[0, 1]$, and any two vertices u, v , define $A_w(u, v)$ to be the probability that u, v are connected in G_w . Here G_w is the subgraph samples from G by picking each edge (u, v) independently with probability $w_{u,v}$.*

Finally, we will also study the popular connectivity measure based on graph effective resistances.

Definition 5. *Let G be a weighted undirected graph with adjacency matrix A and the associated laplacian $L = D - A$. Given an edge u, v the effective resistance between u, v is defined as*

$$R_{u,v} = \xi_{u,v}^\top L^\dagger \xi_{u,v}.$$

Here $\xi_{u,v}$ is an n dimensional vector with $+1$ at position u , -1 at position v and zeros everywhere. L^\dagger refers to the matrix pseudo-inverse.

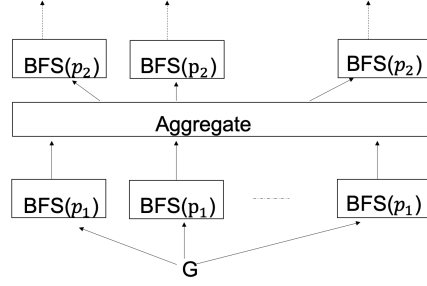
Below we establish the limitations of GNN based architectures in computing the above connectivity measures.

Theorem 10. *Consider a GNN \mathcal{N} of depth L over n nodes where each node has a representation size of B . If for a given pair u, v the GNN \mathcal{N} computes $A(u, v)$, $A_w(u, v)$ or $R_{u,v}$ then it must hold that $B \cdot L \geq \Omega(\frac{\sqrt{n}}{\log n})$.*

Proof. We will establish the above lower bound using tools from communication complexity and in particular the lower bound for the *set disjointness* problem. In the set disjointness problem we are given two inputs $x, y \in \{0, 1\}^b$ each belonging to a different party. The goal is to communicate and determine if $x \cdot y = 0$ or not. It is well known that any protocol that succeeds with constant probability for the set disjointness problem has communication complexity at least $\Omega(b)$ (Babai et al., 1986).

We adapt the construction from Sarma et al. (2012). We construct a graph $G(\Gamma, d, p)$ as follows. There are Γ vertex disjoint paths $\mathcal{P}^1, \mathcal{P}^2, \dots, \mathcal{P}^\Gamma$, each of length d^p . The vertices on path ℓ are denoted by $v_0^\ell, v_1^\ell, \dots, v_{d^p-1}^\ell$. We add to this graph two special vertices s and r . The vertices s and r will be endowed with two Γ dimensional Boolean inputs on which one wants to compute the set disjointness problem. Let $s(i)$ denote the i th input bit for vertex s and $r(i)$ denote the i th input but for vertex r . We connect s to v_0^ℓ for all ℓ such that $s(\ell) = 1$ and r to $v_{d^p-1}^\ell$ for all ℓ such that $r(\ell) = 1$. Furthermore, for a given value of p , we replicated each edge $O(p \log n)$ times.

Next, notice that when we samples edges with probability p , each path edge and edges involving s, r will be sampled at least once with constant probability. Furthermore, conditioned on this event, s and r are connected if and only if the inputs to s and r are not disjoint. Using the simulation theorem of Sarma et al. (2012) we get that for a given value of p , if there is a distributed algorithm that accurately estimates $A_p(u, v)$ with constant probability in $R \leq (d^p - 1)/2$ rounds, then there is a protocol for computing solving set disjointness in at most $O(dpR)$ bits of communication. Hence we

Figure 5: The network architecture for approximate affinity $A(u, v)$.

get that for any p , the number of rounds $(B \cdot L)$ needed to compute $A_p(u, v)$ is at least $\min\left(\frac{d^p-1}{2}, \frac{\Gamma}{dp}\right)$. Writing in terms of $n = \theta(\Gamma d^p)$, and setting $p = O(\log n)$ we get that the number of rounds needed is at least $\Omega\left(\frac{\sqrt{n}}{\log n}\right)$. Since for each p we have a bad input, by minmax we get that the number of rounds $(B \cdot L)$ needed to compute $A(u, v) = \mathbb{E}_p[A_p(u, v)]$ is also at least $\frac{\sqrt{n}}{\log n}$. \square

B.1 UPPER BOUNDS FOR CONNECTIVITY MEASURES

Next we show how to circumvent the above lower bounds via our proposed GNN^+ architecture. We first design an efficient architecture for the connectivity measures $A(u, v)$ and $A_w(u, v)$. For this purpose we adapt the sketch based construction of Panigrahy et al. (2012). The GNN^+ architecture is shown in Figure 5. For an $\epsilon > 0$, the sketch proceeds in $q = O(\log m/\epsilon)$ parallel stages corresponding to the width of the GNN^+ architecture in Figure 5. In each stage q parallel BFS procedures are performed on a sampled graph with edge probability denoted by p_1, p_2, \dots, p_r . At the end of each stage the *Aggregate* procedure maintains a running count of the number of subgraphs in which vertices u and v are in the same connected component. Overall, this is repeated for $r = O(\log m/\epsilon)$ rounds for a total depth of $O(\log m/\epsilon)$. Hence we have the following guarantee.

Theorem 11. *For any constant $\epsilon > 0$, and for any graph over n nodes with maximum degree d and diameter D , the network in Figure 5 produces a $(1 + \epsilon)$ approximation to $A(u, v)$. Furthermore, the network is of size $O(m \log(m))$, has $O(\log m)$ parameters, and depth $O(D \log(m))$.*

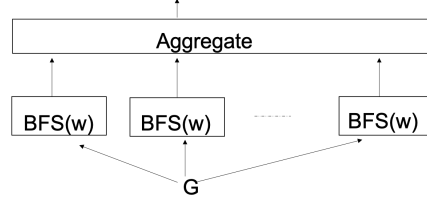
Proof. The proof of correctness of the procedure follows from the analysis of Panigrahy et al. (2012). After each round of parallel BFS procedures, nodes u, v have the id of the connected components that they belong to. The *Aggregate* procedure maintains a running count of in the fraction of connected components in which u and v belonged together. Given $q = O(\log m)$ sized output from the previous stage, this count can be updated by a simple fully connected network with $O(\log m)$ parameters. Furthermore, all the *Aggregate* layers share the same set of parameters. From the previous section, we know that the BFS procedure can be implemented by a GNN of depth $O(D)$ and $O(1)$ total parameters. Hence we get a total depth of $O(D \log m)$ and $O(\log m)$ total parameters. \square

Given the architecture in Figure 5, computing the weighted affinity measure is straightforward by replacing the sampling probability p_1, \dots, p_r with the edge weights. Furthermore, in this case one only needs one stage of parallel BFS procedures. Hence we get the following corollary

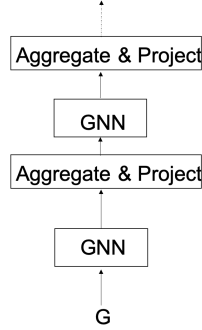
Corollary 2. *For any constant $\epsilon > 0$, and for any graph over n nodes with maximum degree d and diameter D , the network in Figure 6 produces a $(1 + \epsilon)$ approximation to $A_w(u, v)$. Furthermore, the network is of size $O(m \log(m)/w_{\min})$, has $O(\log m)$ parameters, and depth $O(D)$. Here w_{\min} is the minimum edge weight in the graph G .*

Effective Resistance. To approximately compute the effective resistance we consider the following reformulation. Given a graph Laplacian L and vertices u, v consider the following optimization

$$\begin{aligned} \min_{x: \mathbb{R}^n: x \neq 0} x^T L x \\ \text{s.t. } x_u - x_v = 1. \end{aligned} \quad (3)$$

Figure 6: The network architecture for the approximate weighted affinity $A_w(u, v)$.

It is well known that the solution to the above objective equals $1/R_{u,v}$. Hence it is enough to approximate equation 3. Furthermore, equation 3 is a convex optimization problem with a linear constraint and hence can be solved using projected gradient descent. The architecture in Figure 7 implements this via the GNN^+ architecture. In the GNN module each vertex i maintains an $O(d)$

Figure 7: The network architecture for approximate effective resistance $R_{u,v}$.

dimensional representation. At each stage of GNN each vertex computes its local gradient. Notice that given current solution x^t , the gradient of the objective in equation 3 equals Lx^t . Hence, vertex i will compute $L^{(i)}x_i^t$, where $L^{(i)}$ is the i th column of L . The *Aggregate and Project* module first aggregates all the local gradients and performs a gradient step i.e., $x^{t+1} = x^t - \eta Lx^t$ where η is the learning rate. Then it projects x^{t+1} onto the linear subspace $x_u - x_v = 1$. As a result we have the following guarantee

Theorem 12. *For any constant $\epsilon > 0$, and appropriately chosen learning rate η , for any graph over n nodes with maximum degree d and diameter D , the network in Figure 7 produces an additive ϵ approximation to equation 3. Furthermore, the network is of depth $\ell = O(\frac{1}{\lambda_2(L)\epsilon^2} \log n)$, size $O(nd\ell)$, and has $O(d)$ parameters.*

Proof. The number of steps of gradient descent needed for an ϵ approximation equals $\frac{1}{\lambda_2(L)\epsilon^2}$ where $\lambda_2(L)$ bounds the strong convexity of the objective within the linear subspace. In each step, each vertex computes $L^{(i)}x_i^t$ by aggregating messages from its neighbors. This can be done using $O(d)$ parameters. The aggregate step involves first computing $\sum_i L^{(i)}x_i^t$. This can be done by a fully connected network of depth $O(\log n)$ and $O(1)$ parameters. The gradient step can be computed using $O(1)$ parameters depth one. Finally, the project step can be implemented using $O(1)$ parameters since it only involves two coordinates of x^{t+1} . Hence, overall we have $O(d)$ parameters and depth of $O(\frac{1}{\lambda_2(L)\epsilon^2} \log n)$. \square

C CLUSTERING, MIN CUTS AND MST

In this section we demonstrate the power of the GNN^+ architecture for other graph problems such as clustering, cuts and minimum spanning trees. In each case we will achieve significantly better statistical properties than the vanilla GNN^{mp} architectures. We first discuss spectral clustering.

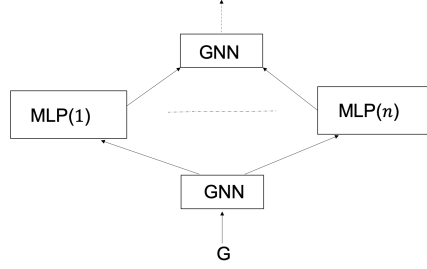


Figure 8: The network architecture for minimum spanning tree (MST).

C.1 SPECTRAL CLUSTERING

. Given an undirected graph G , the spectral clustering of G corresponds to the cut obtained by taking the sign of the eigenvector v corresponding to the second smallest eigenvalue $\lambda_2(L)$. This can be formulated as the following optimization problem

$$\begin{aligned} \min_{x \in \mathbb{R}^n : x \neq 0} x^T L x \\ \text{s.t. } x \perp \mathbf{1}. \end{aligned} \quad (4)$$

From the above formulation it is evident that the GNN^+ based architecture of Figure 7 can be used to efficiently compute an approximate spectral clustering. The only difference between equation 3 and equation 4 is the implementation of the projection step. We discuss this below in the proof of the following theorem

Theorem 13. *For any constant $\epsilon > 0$, and appropriately chosen learning rate η , for any graph over n nodes with maximum degree d and diameter D , the network in Figure 7 produces an additive ϵ approximation to equation 4. Furthermore, the network is of depth $\ell = O(\frac{1}{\lambda_2(L)\epsilon^2} \log n)$, size $O(nd\ell)$, and has $O(d)$ parameters.*

Proof. The proof is exactly the same as that of Theorem 12 the only difference being the implementation of the projection step. In equation 4 given x^{t+1} we need to project it onto the subspace orthogonal to the all ones vector. This corresponds to the following operation

$$x^{t+1} \leftarrow x^{t+1} - \mathbf{1} \frac{x^{t+1} \cdot \mathbf{1}}{n}.$$

The above can be implemented using an $O(\log n)$ depth neural network with $O(1)$ parameters. \square

C.2 MINIMUM SPANNING TREE (MST)

In this section we implement an efficient GNN^+ architecture for computing the minimum spanning tree (MST) in undirected graphs. Note that in the traditional GNN^{mp} framework computing even an approximate MST requires $\Omega(\sqrt{n})$ rounds (Loukas, 2020). Our efficient GNN^+ architecture mirrors the parallel implementation of Boruvka’s algorithm (Sollin, 1965) and is shown in Figure 8. We have the following guarantee associated with the architecture in Figure 8.

Theorem 14. *For any graph over n nodes with maximum degree d and diameter D , the network in Figure 8 produces the minimum weight spanning tree. Furthermore, the network is of depth $\ell = O((D + \log n) \log n)$, size $O(n\ell)$, and has $O(d)$ parameters.*

Proof. The architecture in Figure 8 implements Boruvka’s algorithm (Sollin, 1965) where in each round comprises of a GNN step and an FNN computation. Each vertex maintains an $O(d)$ dimensional representation consisting of the current state of each of its edges. This state can be either *contracted* or *active*. Furthermore, each vertex also maintains a component id that is initialized to be the id of the vertex. Each GNN stage comprises of $O(D)$ rounds of BFS for component id merging. During this, each vertex aggregates messages along contracted edges to update its component id to be the minimum of the component ids received via contracted edges. This can be achieved using $O(D)$

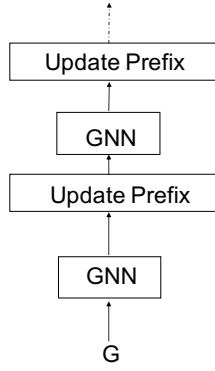


Figure 9: The network architecture for minimum cut.

rounds of GNN computation with $O(d)$ parameters. This is followed by a single round of GNN to figure out the minimum weight edge from each vertex along *active* edges. This information is then fed into the MLP module to decide for each vertex, if any of its active edges will be contracted in the next step. Hence the FNN stage consists of n parallel modules, one per vertex. Each individual MLP module can be implemented using depth $O(\log n)$ and $O(1)$ parameters. Hence, overall we get an architecture of depth $O((D + \log n) \log n)$ and $O(d)$ parameters. \square

C.3 MIN CUT

In this section we design an efficient GNN^+ architecture for computing the global min cut in an undirected graph. Similar to the minimum spanning tree problem, even computing an approximate mincut using traditional GNN^{mp} network requires $\Omega(\sqrt{n})$ depth. As discussed in Section 4, our efficient GNN^+ based architecture is based on the parallel algorithm for computing mincut by Karger & Stein (1996) and is shown in Figure 9. For the proposed architecture we have the following guarantee

Theorem 15 (Restatement of Theorem 9). *For any graph over n nodes with maximum degree d and diameter D , the network in Figure 9 produces the minimum cut. Furthermore, the network is of depth $\ell = O(D \log^2 n)$, size $O(n\ell)$, and has $O(d + \log m)$ parameters.*

Proof. Each vertex maintains an $O(d)$ sized representation indicating which of its edges are currently *active*. The GNN module simply performs a BFS over active edges starting from a specified vertex (output by the *Update Prefix* module). Initially random ordering of the edges is chosen. The goal of the Update Prefix module is to update the current prefix in this ordering that specifies the set of active edges. Let L be the initial ordering of the edges. Then the guarantee from Karger & Stein (1996) implies that with probability at least $1/n^2$ there exists a prefix L' of L such that the subgraph comprising of edges in L' consists of exactly two connected components corresponding to the optimal mincut. Hence, given the output of the BFS performed by the GNN, the Update Prefix module has to decide if the number of connected components is more than two, and update the prefix accordingly. This can be implemented by an FNN using $O(\log n)$ depth and $O(\log m)$ parameters. Hence overall we get $O(D \log n)$ depth and $O(d + \log n)$ parameters. \square

D GENERALIZATION BOUNDS

In this section we derive generalization bounds for the network architectures proposed in this work. In general our network architectures consist of a combination of GNN^{mp} blocks and fully connected layers. We denote by W_i the weight matrix at a fully connected layer and by K_i the kernel matrix of a convolutional operation, either for a GNN or a convolutional layer. Initially, these matrices are initialized to W_i^0 and K_i^0 respectively. Furthermore, for any K_i we will denote by $op(K_i)$ the corresponding linear operator for the kernel matrix K_i . Our generalization bounds will depend on the number of trainable parameters in the network, the depth of the network and the distance of the

parameters to initialization. Let $\mathbf{W}^0, \mathbf{K}^0$ be the initial parameters and let \mathbf{W}, \mathbf{K} be the final trained parameters. We will define the distance between parameters as

$$d(\{\mathbf{W}^0, \mathbf{K}^0\}, \{\mathbf{W}, \mathbf{K}\}) := \sum_i \|W_i^0 - W_i\|_{\infty \rightarrow 2} + \sum_i \|op(K_i^0) - op(K_i)\|_{\infty \rightarrow 2},$$

where for any matrix M , $\|M\|_{\infty \rightarrow 2}$ norm is defined as

$$\|M\|_{\infty \rightarrow 2} := \max_{x: \|x\|_{\infty} \leq 1} \|Mx\|_2.$$

Let $f(G; \mathbf{W}, \mathbf{K})$ be the function, with $\binom{n}{2}$ outputs, defined by a given set of parameters. For any given parameter $\beta > 0$ and initial parameters $\{\mathbf{W}^0, \mathbf{K}^0\}$, we define the class \mathcal{F} to be

$$\mathcal{F}_\beta = \{\{\mathbf{W}, \mathbf{K}\} : d(\{\mathbf{W}^0, \mathbf{K}^0\}, \{\mathbf{W}, \mathbf{K}\}) \leq \beta\}. \quad (5)$$

Our main result will be the following

Theorem 16. *Let $\ell(\hat{y}, y)$ be a λ -Lipschitz loss function bounded in $[0, B]$. Then, given m i.i.d. samples $(G_1, y_1), (G_2, y_2), \dots, (G_m, y_m)$ generated from a distribution D , with probability at least $1 - \delta$ it holds that for all $f \in \mathcal{F}_\beta$,*

$$|\hat{\mathbb{E}}[\ell_f] - \mathbb{E}[\ell_f]| \leq O\left(B\sqrt{\frac{P(\beta + \nu L + \log(\lambda\beta)) + \log(1/\delta)}{m}}\right).$$

Here $1 + \nu$ is the $\|\cdot\|_{\infty \rightarrow 2}$ operator norm bound on the randomly initialized matrices, L is the total depth of the network, and P is the total number of trainable parameters in the network.

As an immediate corollary of the above theorem we get the following guarantees for the shortest path problem.

Corollary 3. *Let \mathcal{F}_β be as defined in equation 5 for the network architecture proposed in Figure 2. Let $\ell(\hat{y}, y)$ be a λ -Lipschitz loss function bounded in $[0, B]$. Then, given m i.i.d. samples $(G_1, y_1), (G_2, y_2), \dots, (G_m, y_m)$ generated from a distribution D , with probability at least $1 - \delta$ it holds that for all $f \in \mathcal{F}_\beta$,*

$$|\hat{\mathbb{E}}[\ell_f] - \mathbb{E}[\ell_f]| \leq O\left(B\sqrt{\frac{n^{2/c}(\beta + \nu(D \log d + \log n) + \log(\lambda\beta)) \log n + \log(1/\delta)}{m}}\right).$$

Corollary 4. *Let \mathcal{F}_β be as defined in equation 5 for the network architecture proposed in Figure 4. Let $\ell(\hat{y}, y)$ be a λ -Lipschitz loss function bounded in $[0, B]$. Then, given m i.i.d. samples $(G_1, y_1), (G_2, y_2), \dots, (G_m, y_m)$ generated from a distribution D , with probability at least $1 - \delta$ it holds that for all $f \in \mathcal{F}_\beta$,*

$$|\hat{\mathbb{E}}[\ell_f] - \mathbb{E}[\ell_f]| \leq O\left(B\sqrt{\frac{n^{2/c}(\beta + \nu(\sqrt{nD} \log d + \log n) + \log(\lambda\beta)) \log n + \log(1/\delta)}{m}}\right).$$

In order to prove Theorem 16 we will use the following general guarantee from the work of Giné & Guillou (2001).

Theorem 17 (Giné & Guillou (2001)). *Let G be a set of functions, parameterized by $\theta \in \mathbb{R}^d$, from a domain Z to $[0, B]$. Furthermore, let G be (γ, d) -Lipschitz, i.e., for some norm $\|\cdot\|$, and for any $g \in G$ and any θ, θ' such that $\|\theta\|, \|\theta'\| \leq 1$, we have*

$$|g_\theta(z) - g_{\theta'}(z)| \leq \gamma \|\theta - \theta'\|$$

. Then for any distribution Z and i.i.d. samples z_1, z_2, \dots, z_m , with probability at least $1 - \delta$, it holds that for any $g \in G$

$$|\hat{\mathbb{E}}[g] - \mathbb{E}[g]| \leq O\left(B\sqrt{\frac{d \log \gamma + \log(1/\delta)}{m}}\right).$$

Using the above theorem we need to show that the class \mathcal{F}_β is Lipschitz. We will establish this next.

Proof of Theorem 16. Consider $f(G; \mathbf{W}, \mathbf{K} \in \mathcal{F}_\beta$ with P total trainable parameters. Hence, the dimensionality of the parameter space is P . We will show that ℓ_f is $(\lambda\beta e^{\nu L + \beta}, P)$ -Lipschitz thereby establishing the theorem statement in light of Theorem 17 above.

Consider parameters \mathbf{W}, \mathbf{K} and \mathbf{W}', \mathbf{K}' . We will move from the two sets of parameters by replacing one parameter at a time. Without loss of generality assume that the parameter K_i at layer i is changed to K'_i . Then we have

$$\begin{aligned} |\ell_f(G; \mathbf{W}, \mathbf{K}) - \ell_f(G; \mathbf{W}', \mathbf{K}')| &\leq \lambda \|f(G; \mathbf{W}, \mathbf{K}) - f(G; \mathbf{W}', \mathbf{K}')\| \\ &\leq \lambda \prod_{j>i} \|W_j\|_{\infty \rightarrow 2} \|op(K_j)\|_{\infty \rightarrow 2} \| (op(K_i) - op(K'_i))u \|. \end{aligned}$$

Here the first inequality follows from the Lipschitz property of the loss function and the second inequality follows from the fact that all layers except i have the same weight parameters. The vector u represents the input to layer i . Since RELU activation is 1-Lipschitz and non-expansive we have

$$\begin{aligned} |\ell_f(G; \mathbf{W}, \mathbf{K}) - \ell_f(G; \mathbf{W}', \mathbf{K}')| &\leq \lambda \prod_{j>i} \|W_j\|_{\infty \rightarrow 2} \prod_{j>i} \|op(K_j)\|_{\infty \rightarrow 2} \|(op(K_i) - op(K'_i))u\| \\ &\leq \lambda \prod_{j \neq i} \|W_j\|_{\infty \rightarrow 2} \prod_{j \neq i} \|op(K_j)\|_{\infty \rightarrow 2} \|(op(K_i) - op(K'_i))\|_{\infty \rightarrow 2} \\ &\leq \lambda \prod_{j \neq i} (\|W_j^0\|_{\infty \rightarrow 2} + \|W_j - W_j^0\|_{\infty \rightarrow 2}) \\ &\quad \times \prod_{j \neq i} (\|op(K_j^0)\|_{\infty \rightarrow 2} + \|op(K_j) - op(K_j^0)\|_{\infty \rightarrow 2}) \|(op(K_i) - op(K'_i))\|_{\infty \rightarrow 2} \\ &\leq \lambda \prod_{j \neq i} (1 + \nu + \gamma_j) \prod_{j \neq i} (1 + \nu + \beta_j) \|(op(K_i) - op(K'_i))\|_{\infty \rightarrow 2}. \end{aligned}$$

Here γ_j is the distance to initialization for W_j and β_j is the distance to initialization for $op(K_j)$. Noting that

$$\sum_j (\gamma_j + \beta_j) = \beta$$

and that the product is maximized when all are equal we get that

$$\begin{aligned} |\ell_f(G; \mathbf{W}, \mathbf{K}) - \ell_f(G; \mathbf{W}', \mathbf{K}')| &\leq \lambda \prod_{j>i} \|W_j\|_{\infty \rightarrow 2} \prod_{j>i} \|op(K_j)\|_{\infty \rightarrow 2} \|(op(K_i) - op(K'_i))u\| \\ &\leq \lambda (1 + \nu + \beta/L)^L \|(op(K_i) - op(K'_i))\|_{\infty \rightarrow 2} \\ &\leq \lambda e^{\nu L + \beta} \|(op(K_i) - op(K'_i))\|_{\infty \rightarrow 2}. \end{aligned}$$

Using triangle inequality and moving one set of parameters at a time we get that ℓ_f is $(\lambda e^{\nu L + \beta}, P)$ -Lipschitz. \square

E ADDITIONAL EXPERIMENTS

In this section, we provide test MSE plots for each of the five graph problems as a function of the number of epochs. While the final test MSE values are presented separately in the main paper and show accuracy gains for the GNN^+ nomdels over GNN^{mp} , these plots here indicate that for most of the problems (except MST), GNN^+ also converges faster than GNN^{mp} .

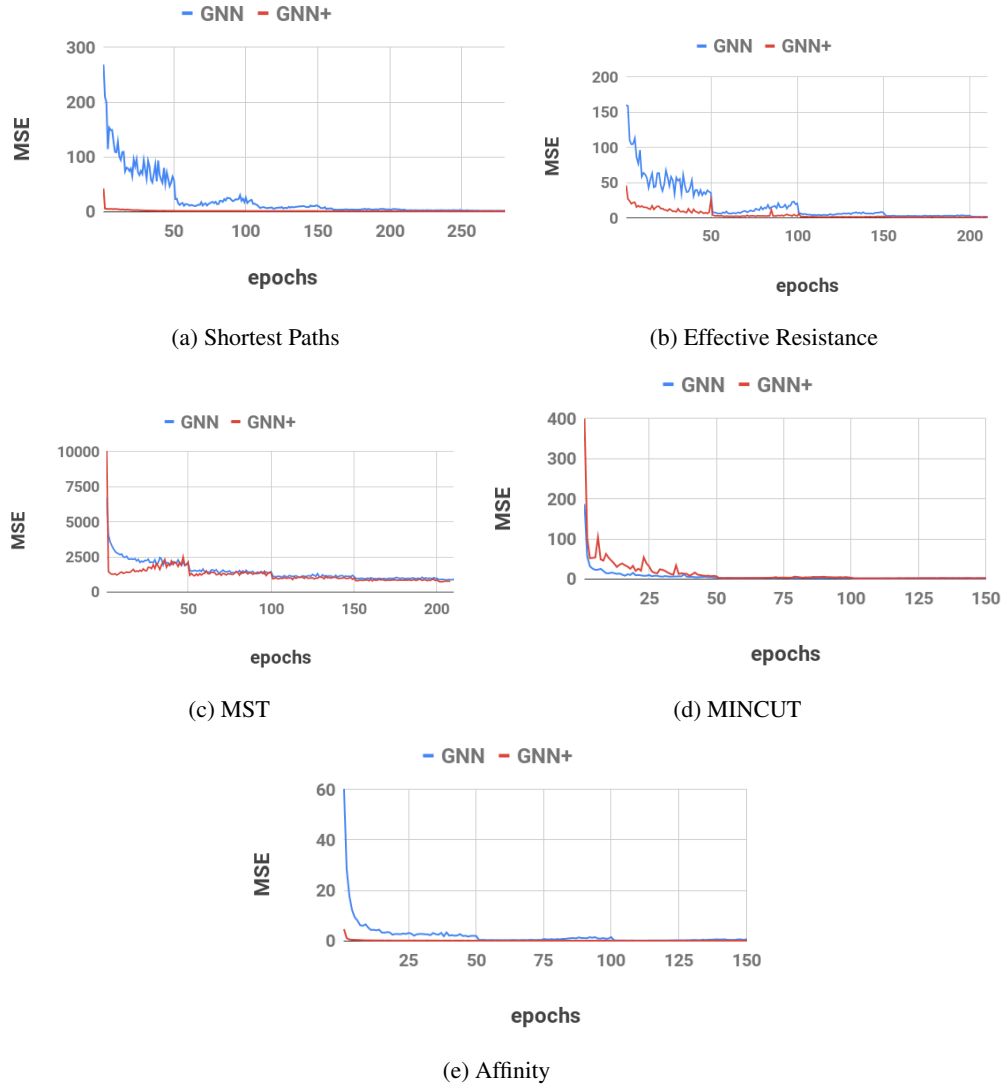


Figure 10: Comparison of GNN^{mp} and GNN^+ on the five graph problems