

# Graph neural networks extrapolate out-of-distribution for shortest paths

**Robert R. Nerem**

RNEREM@UCSD.EDU

*Halcioğlu Data Science Institute, University of California San Diego*

**Samantha Chen**

SAMCHEN@UCSD.EDU

*Department of Computer Science and Engineering, University of California San Diego*

**Sanjoy Dasgupta**

SDASGUPTA@UCSD.EDU

*Halcioğlu Data Science Institute and Department of Computer Science and Engineering, University of California San Diego*

**Yusu Wang**

YUSU@UCSD.EDU

*Halcioğlu Data Science Institute, University of California San Diego*

**Editors:** COLT 2026 Program Chairs

## Abstract

Neural networks (NNs), despite their success and wide adoption, still struggle to extrapolate out-of-distribution (OOD), i.e., to inputs that are not well-represented by their training dataset. Addressing the OOD generalization gap is crucial when models are deployed in environments significantly different from the training set, such as applying Graph Neural Networks (GNNs) trained on small graphs to large, real-world graphs. One promising approach for achieving robust OOD generalization is the framework of *neural algorithmic alignment*, which incorporates ideas from classical algorithms by designing neural architectures that resemble specific algorithmic paradigms (e.g. dynamic programming). The hope is that trained models of this form would have superior OOD capabilities, in much the same way that classical algorithms work for all instances. We employ sparsity regularization as a tool for analyzing the role of algorithmic alignment in achieving OOD generalization, focusing on graph neural networks (GNNs) applied to the canonical *shortest path* problem. We prove that if a trained GNN minimizes a sparsity-regularized loss over a small set of shortest-path instances, then the GNN implements  $K$  steps of the Bellman-Ford algorithm for shortest paths. In fact, if a trained GNN minimizes this loss within an error of  $\epsilon$ , it computes  $K$ -step shortest path distances up to error  $O(\epsilon)$ . Our empirical results support our theory by showing that NNs trained by gradient descent are able to minimize this loss and extrapolate in practice.

**Keywords:** Graph neural networks, algorithmic alignment, out-of-distribution generalization, shortest paths, sparsity regularization

## 1. Introduction

Neural networks (NNs) have demonstrated remarkable versatility across domains, yet a persistent and critical challenge remains in their ability to generalize to out-of-distribution (OOD) inputs, i.e., inputs that differ distributionally from their training data. For instance, a medical diagnosis model trained on North American patients may struggle to generalize when applied to patients in the UK due to differences in underlying population distributions. This issue has motivated entire subfields,

such as distribution shift, transfer learning, and domain adaptation [Hanneke and Kpotufe \(2024\)](#); [Wang et al. \(2021\)](#); [Rosenfeld et al. \(2023\)](#).

Graphs, in particular, highlight this challenge as they can vary dramatically in size, connectivity, and topological features. Graph neural networks (GNNs) [Scarselli et al. \(2008\)](#); [Jegelka \(2022\)](#) have seen tremendous development in the past decade [Zhou et al. \(2020a\)](#); [Wu et al. \(2021\)](#), and have been broadly applied to a wide range of domains, from social network analysis [Hamilton et al. \(2017\)](#); [Fan et al. \(2020\)](#) and molecular property prediction [Gilmer et al. \(2017\)](#); [Wu et al. \(2018\)](#); [Wieder et al. \(2020\)](#) to combinatorial optimization [Cappart et al. \(2023\)](#); [Khalil et al. \(2017\)](#). However, these applications often involve scenarios where the graphs encountered in practice are significantly larger, more complex, or structurally distinct from those in training. The case of **size generalization**, where we hope to generalize to graphs larger than seen in training, is an especially severe case of the OOD generalization problem as the graphs belong to distinct spaces, making the training and test distributions disjoint.

Empirical evidence shows that a powerful route to OOD generalization is *algorithmic alignment*—designing a model’s architecture to match a target algorithmic framework [Xu et al. \(2021\)](#); [Dudzik et al. \(2023\)](#). Such alignment biases the network to finding solutions that resemble algorithms, and thus, inherits properties of algorithms like size independence that can aid OOD generalization. Despite many empirical successes, theoretical guarantees for this approach remain extremely limited.

The primary challenge to OOD guarantees is the very expressivity that makes neural networks successful: highly expressive models can fit the training distribution while realizing hypotheses that fail off-distribution. (See [Yehudai et al. \(2021\)](#) for an example with GNNs, and interesting discussion on the challenges of size generalization.) Existing positive results typically assume a tightly controlled relation between train and test. For example, they embed graphs of different sizes into a shared limit space (graphons) and bound their distance [Ruiz et al. \(2020\)](#); [Levie \(2024\)](#); [Le and Jegelka \(2024\)](#); [Rauchweger and Levie \(2025\)](#), or impose explicit discrepancy bounds between the two distributions [Rosenfeld et al. \(2023\)](#); [Rosenfeld and Garg \(2023\)](#)—thereby narrowing applicability. We instead ask, can OOD generalization guarantees be achieved **without any assumptions on the data**? We answer this question in the affirmative by taking an approach that focuses on analyzing a model’s inductive bias. Unlike PAC learning, where generalization guarantees hold with high probability over a data distribution, our results give deterministic guarantees that hold for all instances. Our results are achieved through explicit regularization that makes the effect of algorithmic alignment on inductive bias analyzable, rendering OOD guarantees tractable.

By combining algorithmic alignment and sparsity (which lend strong inductive biases), we demonstrate that it is possible to train GNNs that provably overcome OOD generalization challenges for the canonical task of computing shortest paths. In particular, we show that if a trained GNN minimizes this loss on just a few well selected small graphs, the resulting GNN implements  $K$  steps of Bellman-Ford and generalizes provably well to larger graphs for  $K$ -step shortest-path computation.

Message-passing graph neural networks are popular architectures for handling data in the form of graphs (cf. surveys [Jegelka \(2022\)](#); [Zhou et al. \(2020b\)](#)). At a high level, they operate by assigning each node  $v$  an *embedding*—say, a vector  $h_v \in \mathbb{R}^d$ —and then iteratively updating these embeddings until they contain a solution to the problem at hand. During each step, each node updates its

embedding based on the embeddings of its neighbors and the weights of the edges connecting them. More precisely, letting  $h_v^{(\ell)}$  denote the embedding after  $\ell$  update steps,

$$h_v^{(\ell)} = f^{\text{up}} \left( h_v^{(\ell-1)}, f^{\text{combine}} \left( \{h_u^{(\ell-1)} \oplus w_{uv} : u \in \mathcal{N}(v)\} \right) \right), \quad (1)$$

where  $\oplus$  denotes concatenation,  $w_{uv}$  is the weight of the edge between  $u$  and  $v$ , the set  $\mathcal{N}(v)$  is the neighborhood of  $v$ , and where  $f^{\text{combine}}$  and  $f^{\text{up}}$  are functions realized by feedforward neural nets ( $f^{\text{combine}}$  is permutation invariant). In this way, after  $\ell$  update steps, each node’s embedding can incorporate information from other nodes up to  $\ell$  hops away (see Figure 1). Since we focus only on message-passing graph neural networks in this work, we refer to them simply as GNNs henceforth. As this model applies to graphs with any number of nodes, a key question is: when and how do GNNs perform well on inputs of varying size?

Neural algorithmic alignment is a well-studied framework aiming to design neural architectures that align structurally with specific algorithmic paradigms for the purpose of improving the OOD generalization abilities of a NN. For instance, an astonishingly vast range of practical algorithms are based on *dynamic programming*, an algorithmic strategy that exploits *self-reducibility* in problems: that is, expressing the solution to the problem in terms of solutions to smaller problems of the same type. Shortest paths admit such a decomposition: if the shortest path from  $s$  to  $t$  goes through node  $u$ , then it consists of the shortest path from  $s$  to  $u$ , followed by the shortest path from  $u$  to  $t$ , two smaller subproblems. Interestingly, dynamic programs appear to be well-aligned with graph neural networks Xu et al. (2020); Dudzik and Veličković (2022).

The Bellman-Ford (BF) algorithm for shortest-path computations is the canonical example in the algorithmic alignment literature, highlighting the connection between dynamic programming and message-passing GNNs. In each iteration  $k$  of BF, the shortest path distances from each node to the source that are achievable with at most  $k$  steps are computed using the shortest path distances achievable with at most  $(k - 1)$  steps. For a specific node  $v$ , the distance  $d_v^{(k)}$  is updated as

$$d_v^{(k)} = \min \left\{ d_u^{(k-1)} + w_{(u,v)} : u \in \mathcal{N}(v) \right\}, \quad (2)$$

where  $w_{(u,v)}$  is the weight of the edge connecting  $u$  to  $v$ . This iterative update process closely mirrors the message-passing mechanism in GNNs, where node features are updated layer by layer based on aggregated information from their neighbors.

**Contribution.** Our work provides theoretical guarantees and empirical validation of out-of-distribution generalization, and marks a significant advancement in understanding the benefits of neural algorithmic alignment. While many prior studies have highlighted the expressivity of NNs, the structural similarity of GNNs and classical algorithmic control flows, and their capacity to mimic algorithmic behavior, they typically fall short of providing rigorous guarantees on generalization, particularly in the OOD setting. In contrast, we show that, by using sparsity regularized training regime, GNNs aligned with the BF algorithm provably extrapolate to arbitrary graphs, regardless of size or structure.

We use sparsity to elucidate the connection between GNNs and the Bellman-Ford (BF) algorithm and demonstrate guaranteed OOD generalization. In particular, we show that if a trained GNN minimizes a sparsity-regularized loss over a particular small set of shortest-path instances, then the GNN implements  $K$  steps of the BF algorithm and hence works on larger graphs for  $K$ -step

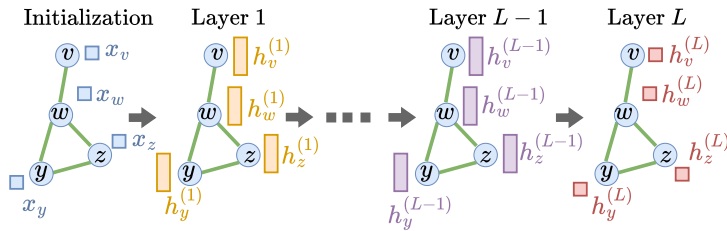


Figure 1: Diagram showing  $L$  GNN layers. The node features are represented by rectangles. At each layer, node features are updated according to neighboring node features and the weights of the adjoining edges.

shortest-path computation. Furthermore, if the trained GNN minimizes the loss up to some error, then it computes  $K$ -step shortest-path distances with at worst proportional error. In this sense, a small sparsity-regularized loss over a specific set of instances serves as a certificate guaranteeing the model’s ability to generalize to arbitrary graphs. Furthermore, we empirically validate that gradient-based optimization indeed finds these BF-aligned solutions, highlighting the practical viability of leveraging algorithmic alignment for enhanced generalization. We believe that our approach combining regularization with algorithmic alignment to analyze OOD generalization will prove useful in other settings. Indeed, subsequent work of [Garnier et al. \(2025\)](#) has used a similar approach to get OOD generalization guarantees for GNNs in solving the heat equation.

### 1.1. Related work

**Neural algorithmic alignment.** Data-dependent approaches to solving combinatorial optimization problems have surged in the past few years [Cappart et al. \(2023\)](#) with GNNs among the most popular architectures used. Early work on GNNs for algorithmic tasks was primarily empirical [Khalil et al. \(2017\)](#); [Karalias and Loukas \(2020\)](#) or focused on representational results [Sato et al. \(2019\)](#); [Loukas \(2020\)](#). The idea of *neural algorithmic alignment* emerged as a conceptual framework for designing suitable GNNs, by selecting architectures that could readily capture classical algorithms for similar tasks. This framework has sample complexity benefits [Xu et al. \(2020\)](#) and promising empirical results [Tang et al. \(2020\)](#); [Grötschla et al. \(2022\)](#); [Xu and Veličković \(2024\)](#). It has also gained traction as a way of understanding the theoretical properties of a given model in terms of its behavior for simple algorithmic tasks (such as BF shortest paths or dynamic programming as a whole) [Veličković et al. \(2022, 2020\)](#); [Veličković and Blundell \(2021\)](#); [Georgiev et al. \(2023\)](#). Our work is the first to establish, both theoretically and empirically, that a NN will converge to the correct parameters which implement a specific algorithm.

**Size generalization.** Size generalization of graph neural networks has been studied empirically, in a variety of settings including classical algorithmic tasks [Veličković et al. \(2022\)](#), physics simulations [Sanchez-Gonzalez et al. \(2020\)](#), and efficient numerical solvers [Luz et al. \(2020\)](#). There is also work on generalization properties of infinite-width GNNs (the so-called *neural tangent kernel* regime); for the simple problem of finding max degree in a graph, [Xu et al. \(2021\)](#) show that graph neural networks in the NTK regime with max readout can generalize to out-of-distribution graphs. Complementary graphon approaches use graph limits combined with continuity of GNNs to understand size generalization [Ruiz et al. \(2020\)](#); [Levie \(2024\)](#); [Le and Jegelka \(2024\)](#); [Rauchwerger](#)

and Levie (2025). Beyond size generalization in GNNs, there is a parallel literature on Transformers that analyzes extrapolation to longer sequences, often called *length generalization* (instead of size generalization). Current work on length generalization asks whether models trained on short sequences of simple arithmetic problems such as addition and modular addition can correctly solve longer problems of the same type Zhou et al. (2023); Jelassi et al. (2023); Liu et al. (2023); Kazemnejad et al. (2024). Within this line, RASP and C-RASP model attention as discrete programs, yielding logic equivalences and depth hierarchies that attempt to explain when such programs extrapolate Weiss et al. (2021); Zhou et al. (2023); Yang and Chiang (2024); Yang et al. (2025). Although these works offer empirically grounded frameworks for size generalization, they provide no formal guarantees. In contrast, we give conditions for provable size generalization.

## 2. Extrapolation Guarantees

### 2.1. Model

In graph neural networks, each node (and possibly edge) is associated with a vector, and in each layer of processing, these vectors are updated based on the vectors of neighboring nodes and adjacent edges. An *attributed* undirected graph is of the form  $G = (V, E, X_v, X_e)$ , where  $X_e = \{x_e : e \in E\}$  are the edge embeddings and  $X_v = \{x_v : v \in V\}$  are the node embeddings. In our case, the edge embeddings will simply be fixed non-negative edge weights,  $x_{(u,v)} = w_{uv}$ , with self-loops set to zero,  $x_{(u,u)} = 0$ . The initial node embeddings  $X_v$  encode the problem input, while the final node embeddings  $h_v^{(L)}$  contain the computed shortest path distances. For instances of shortest path problems, we take  $x_v = 0$  if  $v$  is the source node and use  $x_v = \beta$  to indicate nodes with infinite distance to the source, where  $\beta$  is some number greater than the sum of edge weights. The space of graphs we consider is then

$$\mathcal{G} = \left\{ G = (V, E, X_v, X_e) : \sum_{e \in E} x_e < \beta \right\}.$$

The embedding of node  $v$  at step  $\ell$  is denoted  $h_v^{(\ell)}$  and follows the update rule in Eq. 1 above. (When referring to specific graphs we use  $h_v^{(\ell)}(G)$  and  $x_{(u,v)}(G)$ .)

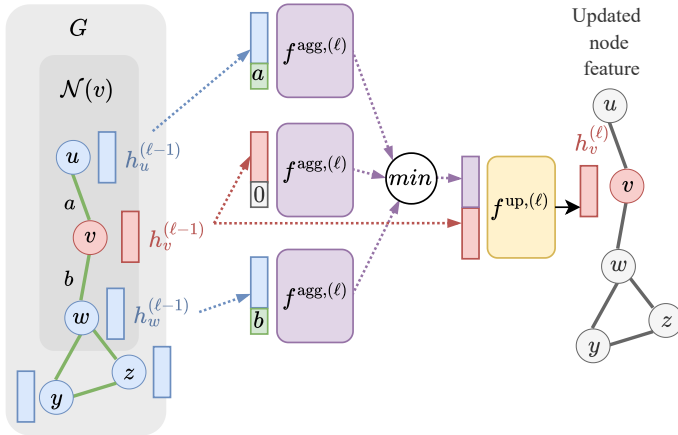


Figure 2: Visual representation of the  $\ell$ -th layer of a MinAgg GNN operating on a graph  $G$ , where  $f^{\text{agg},(\ell)}$  is the aggregation MLP and  $f^{\text{up},(\ell)}$  is the update MLP. Only nodes in the neighborhood  $\mathcal{N}(v)$  of  $v$  are used in the update, so the output at  $v$  is independent of  $x$  and  $y$ .

The  $f^{\text{up},(\ell)}$  function is an MLP that takes two vectors as input: the current embedding of node  $v$ , and a vector representing the aggregated information from  $v$ 's neighbors. It outputs the new embedding of  $v$ . Here  $\mathcal{N}(v)$  denotes the neighbors of node  $v$ , and we take them to include  $v$  itself. The  $f^{\text{combine},(\ell)}$  function combines the embeddings of  $v$ 's neighbors, and the edge weights, into a single vector. A common choice is to apply some MLP  $f^{\text{agg},(\ell)}$  to each (neighbor, edge weight) pair and to then take the sum, or max, or min, of these  $|\mathcal{N}(v)|$  values. We adopt the min. This design choice aligns the network with the structure of the BF algorithm, while still representing a broad and expressive class of GNNs.

**Definition 1** An  $L$ -layer Min-Aggregation Graph Neural Network (MinAgg GNN) with  $d$ -dimensional hidden layers is a map  $\mathcal{A}_\theta : \mathcal{G} \rightarrow \mathcal{G}$  which is computed by layer-wise node-updates (for all  $\ell \in [L]$ ) defined as

$$h_v^{(\ell)} = f^{\text{up},(\ell)} \left( \min_{u \in \mathcal{N}(v)} \{ f^{\text{agg},(\ell)}(h_u^{(\ell-1)} \oplus x_{(u,v)}) \} \oplus h_v^{(\ell-1)} \right) \quad (3)$$

where  $f^{\text{agg},(\ell)} : \mathbb{R}^{d_{\ell-1}+1} \rightarrow \mathbb{R}^d$  and  $f^{\text{up},(\ell)} : \mathbb{R}^{d+d_{\ell-1}} \rightarrow \mathbb{R}^d$  are  $L$ -layer ReLU MLPs, and  $d_0 = d_K = 1$ . Given an input  $G = (V, E, X_v, X_e)$  the initialization is  $h_u^{(0)} = x_u$ .

For simplicity, we assume that  $d_\ell = d$  for  $L > \ell > 0$ . This assumption is made to reduce the number of hyperparameters needed in the analysis, but is made without loss of generality – all of our results hold with general  $d_\ell$ . Furthermore, the choice to make all MLP's have  $L$  layers is also made for simplicity of presentation (again, without loss of generality). Let  $\Gamma$  be a map which implements a single step of the BF algorithm. If  $G = (V, E, X_e, X_v)$  is an attributed graph, then  $\Gamma(G) = (V, E, X_e, X'_v)$  such that for any  $v \in V$ ,

$$x'_v = \min \{ x_u + x_{(u,v)} : u \in \mathcal{N}(v) \cup \{v\} \}.$$

We aim to train a GNN to learn  $K$  iterations of  $\Gamma$ , which we denote by  $\Gamma^K$ .

## 2.2. Warm up

We begin with a warm up that provides intuition for how algorithmic alignment can aid extrapolation. Suppose we look at perhaps the simplest possible GNN that is capable of computing shortest paths. It has updates of the form

$$h_v^{(1)} = \sigma(w_2 \min_{u \in \mathcal{N}(v)} \{ \sigma(W_{11}(x_u \oplus x_{(v,u)} + b_1)) \}) + b_2 \quad (4)$$

Notice that there are just five parameters in this model:  $b_1, b_2, W_{11}, W_{12}$ , and  $w_2$ . The BF algorithm can be simulated by this GNN, as  $b_1 = b_2 = 0$  and  $W_{11} = W_{12} = w_2$  yields Eq. 2. Interestingly, there are many parameter choices that implement the BF update: all that is needed is that  $w_2 W_{11} = w_2 W_{12} = 1$  and  $w_2 b_1 + b_2 = 0$ . Now, let's consider training this model using a small collection of eight graphs, each a path consisting of just one or two edges. Specifically, let  $\mathcal{H}_{\text{small}} = \mathcal{H}_0 \cup \mathcal{H}_1$  with

$$\mathcal{H}_0 = \{ P_1^{(0)}(a_i) : i \in \{1, \dots, 4\} \} \quad \text{and} \quad \mathcal{H}_1 = \{ P_2^{(1)}(a_i, 0) : i \in \{5, \dots, 8\} \} \quad (5)$$

where  $P_k^{(m)}$  denotes a path graph as defined in Fig. 3. The labeled training set is then  $\mathcal{H}_{\text{small}} = \{(G, \Gamma(G)) : G \in \mathcal{H}_{\text{small}}\}$ .

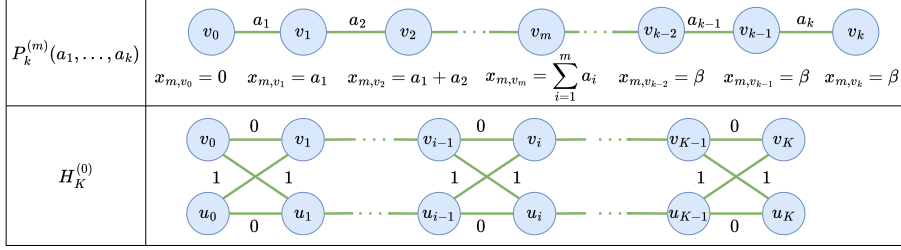


Figure 3: Graphs used in the training sets  $\mathcal{H}_{\text{small}}$  and  $\mathcal{G}_K$ .

**Theorem 2** Let  $0 < \epsilon < 1$ . If  $\forall G \in \mathcal{H}_{\text{small}}$  and  $\forall u \in V(G)$ , a 1-layer GNN  $\mathcal{A}_\theta$  with update given by Eq. 4 computes a node feature satisfying  $|h_u^{(1)}(G) - x_u(\Gamma(G))| < \frac{\epsilon}{20}$ , then for any  $G' \in \mathcal{G}$  and  $v \in V(G')$

$$(1 - \epsilon)x_v(\Gamma(G')) - \epsilon \leq h_v^{(1)}(G') \leq (1 + \epsilon)x_v(\Gamma(G')) + \epsilon.$$

This theorem shows that if the GNN in Eq. 4 achieves low loss on  $\mathcal{H}_{\text{small}}$  then it must implement the  $\Gamma$  operator (a BF step) up to proportionally small error.

**Proof** [Proof Sketch] Note that, if the ReLU activations  $\sigma$  are removed, then eq. (12) is a linear equation. With sufficiently many training examples, this linear equation can fit training data only by implementing the BF algorithm:  $h_v^{(1)} = \min_{u \in \mathcal{N}(v)} (x_u + x_{(v,u)})$ . We proceed by reducing this linear case.

The ReLU activations effectively divides the input space into two halfspaces. This means that the output of the model on any of the input graphs is one of just 4 possible linear functions of the input. The number of input graphs is enough to cover all these cases, so if there is small error on  $\mathcal{H}_{\text{small}}$ , the model must simplify to

$$h_v^{(1)} = w_2 \left( \min_{u \in \mathcal{N}(v)} W_1(x_u \oplus x_{(v,u)} + b_1) + b_2 \right) \quad (6)$$

for most training instances. It is now straightforward to show small error is only achieved if  $w_2 W_{11}$  and  $w_2 W_{12}$  are close to 1 and  $w_2 b_1 + b_2$  is close to zero. These conditions guarantee that the BF algorithm is approximately identified.  $\blacksquare$

### 2.3. Main Result

Now we add one more ingredient, sparsity, which allows us to show extrapolation in *over-parameterized* networks. This time, we consider a full MinAgg GNN as given by Theorem 1. This model is overparameterized in that the MLPs are wider and deeper than needed for BF, and the number of message passing layers is greater than the  $K$  needed to implement  $K$  steps of BF. Consequently, there are a combinatorial number of sparsity patterns, and an infinite number of parameter values, for which this model implements BF, which provides a substantial barrier to achieving and analyzing extrapolation. We introduce sparsity regularization to address this barrier.

To train this model, we again use a small number of simple graphs. The training set contains

$$\mathcal{G}_K = \mathcal{G}_{\text{scale}, K} \cup \{P_1^{(0)}(1), P_2^{(1)}(1, 0), H_K^{(0)}\} \quad (7)$$

where  $\mathcal{G}_{\text{scale}, K}$  contains all path graphs of the form  $P_{K+1}^{(1)}(a, 0, \dots, 0, b, 0, \dots, 0)$  for  $(a, b) \in \{0, 1, \dots, 2K\} \times \{0, 2K + 1\}$  ( $b$  is the weight of the  $k$ th edge). The training instance  $H_K^{(0)}$ , is shown in Fig. 3. The labeled set is  $\mathcal{G}_K = \{(G, \Gamma^K(G)) : G \in \mathcal{G}_K\}$ .

For each graph in the training set  $G \in \mathcal{G}_{\text{train}}$ , we compute the loss only over the set of nodes reachable from the source  $V^*(G)$  (the total number of reachable nodes is  $|\mathcal{G}_{\text{train}}^*|$ ). The regularized loss we use  $\mathcal{L}_{\text{reg}}$  is

$$\mathcal{L}_{\text{reg}}(\mathcal{G}_{\text{train}}, \mathcal{A}_\theta) = \mathcal{L}_{\text{MAE}}(\mathcal{G}_{\text{train}}, \mathcal{A}_\theta) + \eta \|\theta\|_0, \quad (8)$$

where  $\mathcal{L}_{\text{MAE}}$  is

$$\frac{1}{|\mathcal{G}_{\text{train}}^*|} \sum_{G \in \mathcal{G}_{\text{train}}} \sum_{v \in V^*(G)} |x_v(\Gamma^K(G)) - h_v^K(G)|.$$

**Theorem 3** Consider a training set  $\mathcal{G}_{\text{train}}$  with  $M$  total reachable nodes and  $\mathcal{G}_K \subset \mathcal{G}_{\text{train}}$ . For  $L \geq K > 0$ , if an  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$  with  $m$ -layer MLPs achieves a loss  $\mathcal{L}_{\text{reg}}(\mathcal{G}_{\text{train}}, \mathcal{A}_\theta)$  within  $\epsilon$  of its global minimum, where  $0 < \epsilon < \eta < \frac{1}{2M(mL+mK+K)}$ , then on any  $G \in \mathcal{G}$  the features computed by the MinAgg GNN satisfy

$$(1 - M\epsilon)x_v(\Gamma^K(G)) \leq h_v^{(L)}(G) \leq (1 + M\epsilon)x_v(\Gamma^K(G))$$

for all  $v \in V(G)$ .

This theorem shows that low regularized loss implies that an  $L$  layer MinAgg GNN correctly implements  $\Gamma^K$  (i.e.,  $K$ -steps of BF), where the error in implementing this operator is proportional to the distance of the loss from optimal. We later show in experiments that this low loss can be achieved via  $L_1$ -regularized gradient descent. Here we allow for the training set  $\mathcal{G}_{\text{train}}$  to be larger than  $\mathcal{G}_K$  (our guarantee only requires small loss on  $\mathcal{G}_K$ ). However, these additional training examples dilute the training signal from  $\mathcal{G}_K$ , and so the strongest bounds are given if  $\mathcal{G}_{\text{train}} = \mathcal{G}_K$ . We now give a high level road map of our proof, leaving the numerous technical details to Appendix C.

**Proof** [Proof Sketch]

1. **Implementing BF:** We first show  $mL + mK + K$  non-zero parameters are sufficient for the MinAgg GNN to perfectly implement  $K$  steps of the BF algorithm.
2. **Sparsity Constraints:** Next, we show that small loss on  $\mathcal{G}_K$  necessitates a specific sparsity pattern and gives a lower bound of  $mL + mK + K$  on the number of non-zero parameters. Specifically, we show the following.
  - High accuracy on  $P_1^{(0)}(1)$  can only be achieved if each layer of  $f^{\text{up},(\ell)}$  has at least one non-zero entry. *This requires  $mL$  non-zero parameters.*
  - High accuracy on  $H_K^{(0)}$  requires  $K$  layers where  $f^{\text{agg},(\ell)}$  depends on both the node and edge components of its input. This means that each layer of  $f^{\text{agg},(\ell)}$  has at least one non-zero entry, and the first layer of  $f^{\text{agg},(\ell)}$  has two non-zero entries. *This requires  $mK + K$  non-zero parameters.*

We use this fact to derive that the minimum value of  $\mathcal{L}_{\text{reg}}$  is  $\eta(mL + mK + K)$ . The BF implementation reaches the minimum value of  $\mathcal{L}_{\text{reg}}$  since it achieves perfect accuracy with

$mL + mK + K$  non-zero parameters. This result is achieved via a careful selection of training graphs that ensures any model with less than  $mL + mK + K$  non-zero parameters must incur high error on some training set. In particular, we show that for  $K$ -step path graphs

3. **Simplifying the Updates:** Using the above sparsity structure, we can simplify the MinAgg GNN updates to an equivalent update where the intermediate dimensions are always 1 and there are  $K$  updates instead of the previous  $m$  updates:  $\bar{h}_v^{(k)} = \mu^{(k)} \min_{u \in \mathcal{N}(v)} \left\{ \bar{h}_v^{(k-1)} + \nu^{(k)} x_{(u,v)} \right\}$  where  $\mu^{(k)}, \nu^{(k)}, \bar{h}_v^{(k)} \in \mathbb{R}$ . This reduction shows that sparsity regularization addresses the barrier over-parameterization poses to OOD generalization. In particular, simplifying to this linear system allows us to apply the similar techniques as in the warm up to show extrapolation.
4. **Parameter Constraints and Approximation:** If  $\mathcal{L}_{\text{reg}}$  is within  $\epsilon$  of its minimum, the parameters  $\mu^{(k)}, \nu^{(k)}$  must be constrained to avoid poor training accuracy on certain graphs. We then use an induction argument to show these constraints ensure node features approximate BF’s intermediate values and compile these errors to complete the proof. ■

### 3. Experiments

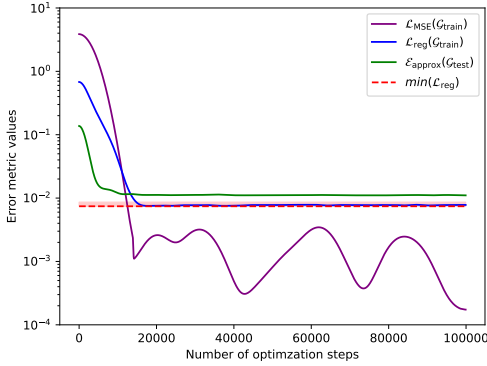
Across all configurations, we use a common training set: the constructed set from our theory (Theorem 3) augmented with a small number of additional graphs to aid optimization in practice. We then vary the OOD test sets, evaluating extrapolation on unseen graphs spanning different sizes and topologies (cycles, complete graphs, and Erdős–Rényi with  $p = 0.5$ ).

Theorems 2 and 3 state that a trained model with a sufficiently low  $L_0$ -regularized loss approximates the BF procedure. We now empirically show that gradient descent can find such a low-loss trained model by optimizing the  $L_1$ -regularized loss  $\mathcal{L}_{\text{MSE}, L_1}$  given by

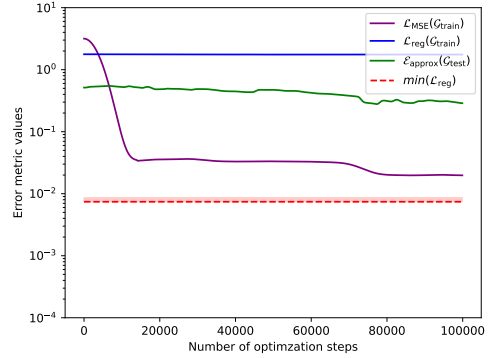
$$\underbrace{\frac{1}{|\mathcal{G}_{\text{train}}^*|} \sum_{G \in \mathcal{G}_{\text{train}}} \sum_{v \in V^*(G)} (x_v(\Gamma^K(G)) - h_v^K(G))^2}_{\mathcal{L}_{\text{MSE}}} + \|\theta\|_1. \quad (9)$$

This training loss  $\mathcal{L}_{\text{MSE}, L_1}$  is a practical proxy for the  $L_0$  regularized loss  $\mathcal{L}_{\text{reg}}$ . To see the effect of sparsity regularization ( $L_1$ -term), we also train a comparison model using the *unregularized loss*  $\mathcal{L}_{\text{MSE}}$  (bracketed terms in Eq. (9)). We show that models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  find sparse generalizable solutions; while models trained without sparsity regularization have worse generalization.

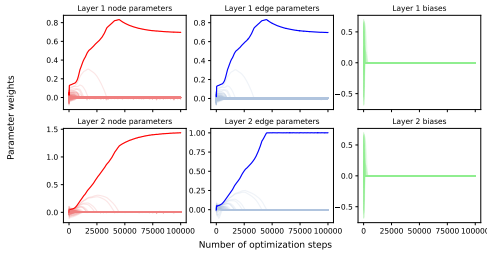
**Additional setup.** We verify our theoretical results empirically using a two-layer MinAgg GNN trained on two steps of BF. We focus on two steps of BF in the main experiments for computational tractability, while using repeated applications of the trained model to evaluate behavior beyond two steps. Specifically, we show that converging to a low value of  $\mathcal{L}_{\text{reg}}$  indicates better performance – particularly in improving generalization to larger test graphs. We additionally show that with  $L_1$  regularization, the trained model parameters approximately implement a sparse BF step. In our experiments, we configure the MinAgg GNN with two layers and 64 hidden units in both the aggregation and update functions. The first layer has an output dimension of eight, while the



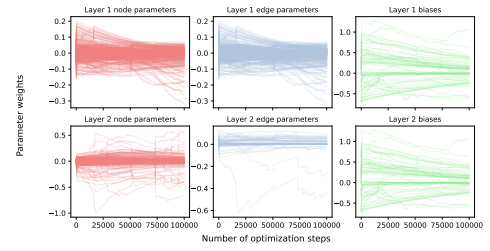
(a) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .



(b) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}}$ .



(c) Model parameters summaries for model trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .



(d) Model parameters summaries for model trained with  $\mathcal{L}_{\text{MSE}}$ .

Figure 4: Performance metrics and parameter updates for a two-layer MinAgg GNN trained on two steps of the BF algorithm. The dotted line in (a) and (b) is the global minimum of Eq. (8). In (a) and (b), we track the change in the train loss, test loss, and  $\mathcal{L}_{\text{reg}}$  over each optimization step for the models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ . (c) and (d) show changes in model parameters over each optimization step with and without  $L_1$  regularization, respectively. Each curve has been smoothed with a truncated Gaussian filter with  $\sigma = 20$ .

second layer outputs a single value. In the supplement, we present additional results evaluating the performance of several other model configurations on one and two steps of BF. To evaluate trained models, we use the following three error metrics:

1. **Empirical training error** ( $\mathcal{L}_{\text{MSE}}$ ): We track the training loss for the model on  $\mathcal{G}_{\text{train}}$ .  $\mathcal{G}_{\text{train}}$  consists of  $\mathcal{G}_K$  where  $K = 2$  as well as four three-node path graphs initialized at step zero of BF and four five-node path graphs initialized at step two of BF. We include these extra graphs to provide examples for the initial and final two steps of the BF algorithm. Empirically, we observe that this expanded training set eases model convergence.
2. **Test error** ( $\mathcal{E}_{\text{approx}}$ ): We compute the average relative error of the model predictions compared to the ground-truth BF output over a test set  $\mathcal{G}_{\text{test}}$ :

$$\mathcal{E}_{\text{approx}}(\mathcal{G}_{\text{test}}) = \frac{1}{|\mathcal{G}_{\text{test}}|} \sum_{G \in \mathcal{G}_{\text{test}}} \sum_{v \in V(G)} \left| \frac{h_v^K - x_v(\Gamma^K(G))}{h_v^K(G)} \right|.$$

# of nodes	Single		Iterated	
	No $L_1$ -reg.	With $L_1$ -reg.	No $L_1$ -reg.	With $L_1$ -reg.
100	0.0202 $\pm$ 0.0055	<b>0.0014 <math>\pm</math> 0.0002</b>	0.0617 $\pm$ 0.0111	<b>0.0036 <math>\pm</math> 0.0002</b>
500	0.0242 $\pm$ 0.0231	<b>0.0015 <math>\pm</math> 0.0002</b>	0.0881 $\pm$ 0.0080	<b>0.0036 <math>\pm</math> 0.0002</b>
1K	0.0183 $\pm$ 0.0066	<b>0.0035 <math>\pm</math> 0.0005</b>	0.0951 $\pm$ 0.0127	<b>0.0092 <math>\pm</math> 0.0030</b>

Table 1: Measuring  $\mathcal{E}_{\text{approx}}$  as the number of nodes per graph increases. We test models trained with  $\mathcal{L}_{\text{MSE},L_1}$  and models trained with  $\mathcal{L}_{\text{MSE}}$ . For each model, we examine  $\mathcal{E}_{\text{approx}}$  (first two columns): for two steps of BF (a single forward pass of each model) and (last two columns): for six steps of BF (where each model is iterated three times). Each test set consists of Erdős–Rényi graphs generated with the corresponding sizes listed with  $p$  such that the expected degree  $np = 5$ . For both models, there is little variation in  $\mathcal{E}_{\text{approx}}$  as the graph size increases. However, for the iterated version of each model,  $\mathcal{E}_{\text{approx}}$  for the model trained with  $\mathcal{L}_{\text{MSE},L_1}$  remains accurate, while the unregularized model shows a significantly larger test error when iterated 3 times.

$\mathcal{G}_{\text{test}}$  consists 200 total graphs. In order to test the generalization ability of each model, we construct  $\mathcal{G}_{\text{test}}$  from 3-cycles, 4-cycles, complete graphs (with up to 200 nodes), and Erdős–Rényi graphs generated using  $p = 0.5$ .

3.  **$L_0$ -regularized error ( $\mathcal{L}_{\text{reg}}$ ):** This metric, which is  $\mathcal{L}_{\text{reg}}$  (see Eq. 8) evaluated on  $\mathcal{G}_{\text{train}}$ , shows how the model’s performance satisfies the conditions of Theorem 3.

We also show model parameters per epoch. In brief, at each layer, we track biases, the parameters that scale the node features, and the parameters that scale the edge features. For the sparse implementation of BF, the node and edge parameter updates both have the same single non-zero positive value  $a$  in the first layer. In the second layer, the node and edge parameter updates both have a single non-zero positive value but the edge parameter update converges to 1 while the node parameter update converges to  $1/a$ .

**Results.** Fig. 4 shows the results of training on two steps of BF. Here (a) and (b) show  $\mathcal{L}_{\text{MSE},L_1}$  (i.e., the model trained with regularized loss) achieves a low value of  $\mathcal{L}_{\text{reg}}$  and a correspondingly a low test error,  $\mathcal{E}_{\text{approx}}$  (in the supplement, we show that this small value of  $\mathcal{L}_{\text{reg}}$  satisfies the conditions of Theorem 3). In contrast, the model trained with  $\mathcal{L}_{\text{MSE}}$  (i.e., unregularized loss) has significantly higher  $\mathcal{L}_{\text{reg}}$  and  $\mathcal{E}_{\text{approx}}$ . Although both models achieve values of  $\mathcal{L}_{\text{MSE}}$  below 0.10, we see that the train error does not necessarily indicate if a sparse implementation of Bellman-Ford has been learned. As such, Fig. 4 (a) and (b) experimentally validates Theorem 3 by demonstrating that low values of  $\mathcal{L}_{\text{reg}}$  yield better generalization on test graphs with different sizes and topologies from the train graphs. In Fig. 4 (c) and (d), we elucidate the effect of  $L_1$  regularization with regards to achieving low values of  $\mathcal{L}_{\text{reg}}$  and show that the model trained with  $\mathcal{L}_{\text{MSE},L_1}$  indeed approximates a sparse implementation of BF. For these sparsity measurements, weights with magnitude below  $10^{-4}$  are rounded to zero when computing the experimental  $\ell_0$  norm. To further illustrate that achieving low values of  $\mathcal{L}_{\text{reg}}$  implies that the BF GNN model will learn to implement Bellman-Ford, we provide parameter heatmaps for models trained with  $L_1$  regularization and without in Fig. 5. We observe that the model trained with  $L_1$  regularization implements exactly the parameters for Bellman-Ford suggested from our theory.

In Table 1, we further assess the generalization ability of the  $L_1$ -regularized model on sparse Erdős-Renyí graphs of increasing size. When we use the trained 2-step MinAgg GNN as a primitive module and iterate it 3 times (to estimate  $2 \times 3 = 6$  BF steps), the test error of our  $L_1$ -regularized model *does not accumulate* while the error by the un-regularized model increases by roughly a factor of 3. Again,  $L_1$  regularization improves generalization. By iteratively applying the trained 2-step BF multiple times, we obtain a neural model to approximate general shortest paths with guarantees.

#### 4. Discussion

We show that algorithmic alignment can fundamentally enhance out-of-distribution (OOD) generalization. By training GNNs with a sparsity-regularized loss on a small set of shortest-path instances, we obtain models that correctly implement the BF algorithm. This result provides a theoretical guarantee that the learned network can generalize OOD to graphs of sizes and structures beyond those encountered during training. This is one of the first results where a neural model, when trained to sufficiently low loss, can **guarantee OOD size generalization** for a non-linear algorithm.

A key challenge in machine learning research lies in comparing neural networks with similar expressivity but different generalization behaviors. This challenge centers on understanding the inductive biases that guide models toward particular solutions. Sparsity regularization creates a setting where the influence of architecture choice on inductive bias is clear and easy to analyze. Indeed, our work demonstrates how this regularization interacts with the GNN architecture to bias the GNN toward implementing the Bellman-Ford algorithm. By making inductive biases explicit and quantifiable, we gain insight into why algorithmic alignment is effective at promoting neural networks to generalize beyond their training distribution, a critical capability for real-world applications.

Our theoretical analysis uses  $\ell_0$  regularization because it gives a clean way to relate minimizers of the loss to sparse parameter structures that implement Bellman-Ford. In the experiments, we use  $L_1$  regularization as a simple practical proxy for inducing this sparsity, and observe that it drives the model toward the sparse BF implementation identified by the theory. Characterizing the minimizers of the corresponding  $L_1$ -regularized objective remains an important open direction. More broadly, one could also consider other sparsity-promoting methods for neural networks, including pruning-based approaches, variational dropout, and relaxations or reparameterizations of  $\ell_0$  regularization [Han et al. \(2015\)](#); [Molchanov et al. \(2017\)](#); [Louizos et al. \(2018\)](#); [Li and Ji \(2019\)](#). Recent work on  $\ell_p$  regulariza-

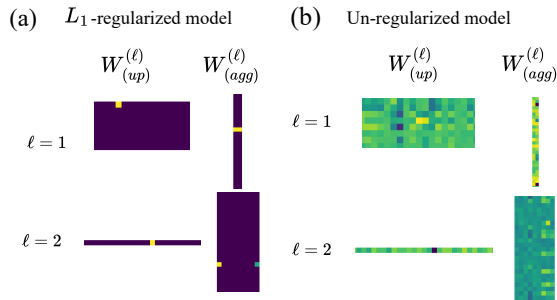


Figure 5: Example of the parameter heatmaps for both the  $L_1$ -regularized model and the un-regularized model. Our Bellman-Ford model update per node feature is defined as  $\sigma(W_{(up)}^{(\ell)} \min\{\sigma(W_{(agg)}^{(\ell)}(x_u + x_{(u,v)} + b_{(agg)}^{(\ell)}) : u \in \mathcal{N}(v)\} + b_{(up)}^{(\ell)}$ . Notice that our  $L_1$  regularized model only has 6 non-zero parameters (exactly as suggested by our theoretical results) while the un-regularized model does not.

tion with  $p < 1$  suggests another promising direction, showing in simple ReLU settings that such objectives can recover sparsest minimizers [Nakhleh and Nowak \(2025\)](#).

**Extensions.** As our BF-aligned GNN correctly implements  $K$  steps, we can extend this capability by recurrently iterating the network. (See the last column of Table 1.) This approach allows the network to solve shortest-path problems that require more than  $K$  iterations. Such scalability enables generalization to shortest-path computations that require arbitrary computational costs.

Furthermore, the ability to learn a single algorithmic step is valuable in broader contexts of neural algorithmic reasoning. This modular design means that the MinAgg GNN can serve as a subroutine within more complex neural architectures that aim to solve higher-level tasks. For instance, in neural combinatorial optimization or graph-based decision-making tasks, shortest-path computations are often just one component of a larger process. By ensuring the network reliably implements each step of the BF algorithm, we create a reusable building block that can be integrated into more sophisticated models. This supports the goal of developing NNs that can reason algorithmically, enabling them to solve increasingly complex problems through the composition of learned algorithmic steps.

**Conclusion.** Our work opens an exciting new direction for research by raising the question of when low training loss can serve as a guarantee for out-of-distribution generalization in other tasks or architectures. While our results focus on the BF algorithm and message-passing GNNs, they suggest the potential for similar guarantees in other instances of alignment, such as different (dynamic programming based) graph algorithms, sequence-to-sequence tasks, or architectures like transformers and recurrent neural networks. Investigating the structural and algorithmic properties that enable such guarantees could provide a unified framework for designing NNs that generalize reliably across diverse tasks and input domains.

## Acknowledgements

This work was supported by the US National Science Foundation under grants CCF-2112665, CCF-2217058, and CCF-2310411.

## References

- Quentin Cappart, Didier Chételat, Elias B Khalil, Andrea Lodi, Christopher Morris, and Petar Veličković. Combinatorial optimization and reasoning with graph neural networks. *Journal of Machine Learning Research*, 24(130):1–61, 2023.
- Andrew Dudzik, Tamara von Glehn, Razvan Pascanu, and Petar Velivckovi’c. Asynchronous algorithmic alignment with cocycles. *ArXiv*, abs/2306.15632, 2023. URL <https://api.semanticscholar.org/CorpusID:259262557>.
- Andrew J Dudzik and Petar Veličković. Graph neural networks are dynamic programmers. In *Advances in Neural Information Processing Systems*, volume 35, pages 20635–20647, 2022.
- Wenqi Fan, Yao Ma, Qing Li, Jianping Wang, Guoyong Cai, Jiliang Tang, and Dawei Yin. A graph neural network framework for social recommendations. *IEEE Transactions on Knowledge and Data Engineering*, 34(5):2033–2047, 2020.

- Paul Garnier, Jonathan Viquerat, and Elie Hachem. Automated discovery of finite volume schemes using graph neural networks, August 2025. URL <https://arxiv.org/abs/2508.19052>.
- Dobrik Georgiev Georgiev, Pietro Lio, Jakub Bachurski, Junhua Chen, Tunan Shi, and Lorenzo Giusti. Beyond erdos-renyi: Generalization in algorithmic reasoning on graphs. In *The Second Learning on Graphs Conference*, 2023.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272. PMLR, 2017.
- Florian Grötschla, Joël Mathys, and Roger Wattenhofer. Learning graph algorithms with recurrent graph neural networks. *arXiv preprint arXiv:2212.04934*, 2022.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. In *Advances in Neural Information Processing Systems*, volume 28, 2015.
- Steve Hanneke and Samory Kpotufe. A more unified theory of transfer learning. *arXiv preprint arXiv:2408.16189*, 2024.
- Stefanie Jegelka. Theory of graph neural networks: representation and learning. In *The International Congress of Mathematicians*, pages 1–23, 2022.
- Samy Jelassi, Stéphane d’Ascoli, Carles Domingo-Enrich, Yuhuai Wu, Yuanzhi Li, and François Charton. Length generalization in arithmetic transformers. *arXiv preprint arXiv:2306.15400*, 2023.
- Nikolaos Karalias and Andreas Loukas. Erdos goes neural: an unsupervised learning framework for combinatorial optimization on graphs. In *Advances in Neural Information Processing Systems*, volume 33, pages 6659–6672, 2020.
- Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. In *Advances in Neural Information Processing Systems*, volume 36, 2024.
- Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial optimization algorithms over graphs. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Thien Le and Stefanie Jegelka. Limits, approximation and size transferability for gnns on sparse graphs via graphops. In *Advances in Neural Information Processing Systems*, volume 36, 2024.
- Ron Levie. A graphon-signal analysis of graph neural networks. In *Advances in Neural Information Processing Systems*, volume 36, 2024.

- Yang Li and Shuiwang Ji.  $L_0$ -arm: Network sparsification via stochastic binary optimization. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, Lecture Notes in Computer Science. Springer, 2019.
- Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=De4FYqjFueZ>.
- Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through  $l_0$  regularization. In *International Conference on Learning Representations*, 2018.
- Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International Conference on Learning Representations*, 2020.
- Ilay Luz, Meirav Galun, Haggai Maron, Ronen Basri, and Irad Yavneh. Learning algebraic multigrid using graph neural networks. In *International Conference on Machine Learning*, pages 6489–6499. PMLR, 2020.
- Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pages 2218–2227, 2017.
- John B. Nakhleh and Robert D. Nowak. Global minimizers of  $\ell_p$ -regularized objectives yield the sparsest ReLU neural networks. In *Advances in Neural Information Processing Systems*, 2025. Poster.
- Levi Rauchwerger and Ron Levie. A note on graphon-signal analysis of graph neural networks. *arXiv preprint arXiv:2508.18564*, 2025. URL <https://arxiv.org/abs/2508.18564>.
- Elan Rosenfeld and Saurabh Garg. (almost) provable error bounds under distribution shift via disagreement discrepancy. In *Advances in Neural Information Processing Systems*, volume 36, 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/5bacb12bf81e98e2ee0eed953a23c656-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/5bacb12bf81e98e2ee0eed953a23c656-Paper-Conference.pdf).
- Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. Domain-adjusted regression or: Erm may already learn features sufficient for out-of-distribution generalization. In *International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=4XE614GBuGR>.
- Ryan A. Rossi and Nesreen K. Ahmed. Usair97 (u.s. air transportation network, 1997). <http://networkrepository.com/USAir97.php>. Accessed 2025-11-24.
- Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2015. URL <http://networkrepository.com>.
- Luana Ruiz, Luiz F. O. Chamon, and Alejandro Ribeiro. Graphon neural networks and the transferability of graph neural networks. In *Advances in Neural Information Processing Systems*, volume 33, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/12bcd658ef0a540cab36cdf2b1046fd-Paper.pdf>.

- Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and Peter Battaglia. Learning to simulate complex physics with graph networks. In *International Conference on Machine Learning*, pages 8459–8468. PMLR, 2020.
- Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Approximation ratios of graph neural networks for combinatorial problems. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.
- Hao Tang, Zhiao Huang, Jiayuan Gu, Bao-Liang Lu, and Hao Su. Towards scale-invariant graph-related problem solving by iterative homogeneous gnns. In *Advances in Neural Information Processing Systems*, volume 33, pages 15811–15822, 2020.
- Petar Veličković and Charles Blundell. Neural algorithmic reasoning. *Patterns*, 2(7), 2021.
- Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural execution of graph algorithms. In *International Conference on Learning Representations*, 2020.
- Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino, Misha Dashevskiy, Raia Hadsell, and Charles Blundell. The clrs algorithmic reasoning benchmark. In *International Conference on Machine Learning*, pages 22084–22102. PMLR, 2022.
- Tianjin Wang, Yongchao Zhang, Renjie Liao, Lawrence Carin, Zheyang Chen, and Andrew G Wilson. Towards out-of-distribution generalization: a survey. *arXiv preprint arXiv:2108.13624*, 2021.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. *arXiv preprint arXiv:2106.06981*, 2021. URL <https://arxiv.org/abs/2106.06981>.
- Oliver Wieder, Stefan Kohlbacher, Méline Kuenemann, Arthur Garon, Pierre Ducrot, Thomas Seidel, and Thierry Langer. A compact review of molecular property prediction with graph neural networks. *Drug Discovery Today: Technologies*, 37:1–12, 2020.
- Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical Science*, 9(2):513–530, 2018.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2021. doi: 10.1109/TNNLS.2020.2978386.
- Kaijia Xu and Petar Veličković. Recurrent aggregators in neural algorithmic reasoning. *arXiv preprint arXiv:2409.07154*, 2024.
- Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? In *International Conference on Learning Representations*, 2020.

Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. How neural networks extrapolate: From feedforward to graph neural networks. In *International Conference on Learning Representations*, 2021.

Andy Yang and David Chiang. Counting like transformers: Compiling temporal counting logic into softmax transformers. *arXiv preprint arXiv:2404.04393*, 2024. URL <https://arxiv.org/abs/2404.04393>.

Andy Yang, Guillaume Cadilhac, and David Chiang. Knee-deep in c-rasp: A transformer depth hierarchy. *arXiv preprint arXiv:2504.12345*, 2025. URL <https://arxiv.org/abs/2504.12345>.

Gilad Yehudai, Ethan Fetaya, Eli Meirom, Gal Chechik, and Haggai Maron. From local structures to size generalization in graph neural networks. In *International Conference on Machine Learning*, pages 11975–11986. PMLR, 2021.

Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio, and Preetum Nakkiran. What algorithms can transformers learn? a study in length generalization. *arXiv preprint arXiv:2310.16028*, 2023.

Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020a. doi: 10.1016/j.aiopen.2021.01.001.

Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: a review of methods and applications. *AI Open*, 1:57–81, 2020b. ISSN 2666-6510. doi: <https://doi.org/10.1016/j.aiopen.2021.01.001>. URL <https://www.sciencedirect.com/science/article/pii/S2666651021000012>.

## Appendix A. Definitions and notation

We begin with definitions and notations we utilize in our proof. To ensure the supplementary material is easy to navigate and reader-friendly, we reiterate some definitions from the main text. We take  $[n] = \{1, 2, \dots, n\}$  and use  $x \oplus y$  to denote the concatenation of the vectors  $x$  and  $y$ . The neighborhood of a node  $v$  is denoted  $\mathcal{N}(v)$  and we use the convention  $v \in \mathcal{N}(v)$ . When referring to the  $i$ th component of  $x$  we write  $x_i$  or  $[x]_i$ .

Given a source node  $s \in V$ , let  $d^{(t)}(s, v)$  denote the length of the  $t$ -step shortest path from  $s$  to  $v$ . If no such path exists,  $d^{(t)}(s, v) = \beta$  and  $\beta$  is some large number. We define a single  $t$ -step Bellman-Ford instance to be a attributed graph  $G^{(t)} = (V, E, X_v, X_e)$  where  $X_v = \{x_v = d^{(t)}(s, v) : v \in V\}$  for some  $s \in V$ . For every 0-step Bellman-Ford instance  $G^{(0)} = (V, E, X_v, X_e)$ ,  $x_s = 0$  for the source node  $s \in V$  and  $x_u = \beta$  for all other nodes  $u \in V$ . Throughout this manuscript,  $t$ -step BF instances are always denoted by a superscript  $(t)$ . Recall that all edge weights considered in this manuscript are non-negative.

Let  $\Gamma$  be a map which implements a single step of the BF algorithm. If  $G = (V, E, X_e, X_v)$  is an attributed graph, then  $\Gamma(G) = (V, E, X_e, X'_v)$  such that for any  $v \in V$ ,

$$x'_v = \min\{x_u + x_{(u,v)} : u \in \mathcal{N}(v)\}.$$

Let  $\Gamma^K$  be  $K$  iterations of  $\Gamma$ . Note that applying  $\Gamma^K$  to a 0-step Bellman-Ford instance  $G^{(0)}$  yields the  $K$ -step shortest path from  $s$  to  $v$ , i.e.,  $\Gamma^K(G^{(0)}) = G^{(K)}$ . Although we restrict our training set to BF instances, our extrapolation guarantees show that the MinAgg GNN approximates the operator  $\Gamma^K$  on any graph in

$$\mathcal{G} = \left\{ G = (V, E, X_v, X_e) : \sum_{e \in E} x_e < \beta \right\}.$$

Define a length- $k$  path graph instance as  $P_k^{(t)}(a_1, \dots, a_k) = (V, E, X_v, X_e)$  where  $V = \{v_0, v_1, \dots, v_k\}$  and  $E = \{(v_{i-1}, v_i) \mid i \in \{1, \dots, k\}\}$ . Let  $x_{(v_{i-1}, v_i)} = a_i$ ,  $x_{v_0} = 0$  (i.e. the source node is  $s = v_0$ ) and  $x_{v_i} = d^{(t)}(s, v_i)$  for  $i > 0$ .

**Definition 4** An  $L$ -layer MinAgg GNN with  $d$ -dimensional hidden layers is a map  $\mathcal{A}_\theta : \mathcal{G} \rightarrow \mathcal{G}$  which is computed by layer-wise node-updates (for all  $\ell \in [L]$ ) defined as

$$h_v^{(\ell)} = f^{\text{up},(\ell)} \left( \min_{u \in \mathcal{N}(v)} \{f^{\text{agg},(\ell)}(h_u^{(\ell-1)} \oplus x_{(u,v)})\} \oplus h_v^{(\ell-1)} \right) \quad (10)$$

where  $f^{\text{agg},(\ell)} : \mathbb{R}^{d_{\ell-1}+1} \rightarrow \mathbb{R}^d$  and  $f^{\text{up},(\ell)} : \mathbb{R}^{d+d_{\ell-1}} \rightarrow \mathbb{R}^{d_\ell}$  are  $L$ -layer ReLU MLPs, and  $d_0 = d_K = 1$ . Given an input  $G = (V, E, X_v, X_e)$  the initialization is  $h_u^{(0)} = x_u$ . The MinAgg GNN  $\mathcal{A}_\theta$  has output  $\mathcal{A}_\theta(G) = (V, E, X'_v = \{h_v^{(\ell)} : v \in V\}, X_e)$ . A **simple**  $L$ -layer MinAgg GNN instead uses the layer-wise update

$$h_v^{(\ell)} = f^{\text{up},(\ell)} \left( \min_{u \in \mathcal{N}(v)} \{f^{\text{agg},(\ell)}(h_u^{(\ell-1)} \oplus x_{(u,v)})\} \right). \quad (11)$$

We refer to the first  $d_{\ell-1}$  components of the domain of  $f^{\text{agg},(\ell)}$  as its node component, and we refer to the last component of the domain of  $f^{\text{agg},(\ell)}$  as its edge component.

A  $K$ -step training set  $\mathcal{G}_{\text{train}}$  is a set of tuples where for each element  $(G^{(t)}, \Gamma^K(G^{(t)})) \in \mathcal{G}_{\text{train}}$  the graph  $G^{(t)}$  is a  $t$ -step BF instance. For a graph  $G = (V, E, X_v, X_e)$  let  $V^*(G) = \{v \in V : x_v \neq \beta\}$  be the set of reachable nodes and let  $|\mathcal{G}_{\text{train}}|^* = \sum_{G^{(t)} \in \mathcal{G}_{\text{train}}} |V^*(G^{(t)})|$  be the total number of reachable nodes in the training set. For each graph  $G$  we consider the training loss over the subset of vertices  $V^*(G)$  because the choice of the feature at unreachable nodes  $\beta$  is arbitrary and so should not be included when providing supervision for shortest path problems.

**Definition 5** An  $m$ -layer ReLU MLP is a function  $f_\theta : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_m}$  parameterized by  $\theta = \{W_j : W_j \in \mathbb{R}^{d_j \times d_{j-1}}, j \in [m]\} \cup \{b_j : b_j \in \mathbb{R}^{d_j}, j \in [m]\}$  where for all  $j \in [m]$ ,

$$\begin{aligned} x^{(0)} &= x, \\ x^{(j)} &= \sigma(W_j x^{(j-1)} + b_j), \end{aligned}$$

and  $f_\theta(x) = x^{(m)}$ . Here,  $\sigma$  is the rectified linear unit (ReLU) activation function.

The MinAgg GNN is parameterized by the set of weights

$$\theta = \bigcup_{\ell=1}^L (\theta^{\text{up},(\ell)} \cup \theta^{\text{agg},(\ell)}),$$

where  $\theta^{\text{up},(\ell)}$  and  $\theta^{\text{agg},(\ell)}$  denote the parameters of the update and aggregation MLPs at layer  $\ell$ , respectively.

We also formalize the definition of path graph instances. A 0-step path graph instance  $P_k^{(0)}(a_1, \dots, a_k)$  consists of a graph  $(V, E, X_v, X_e)$  where the vertex set is  $V = \{v_0, v_1, \dots, v_k\}$ , the edge set is  $E = \{(v_{i-1}, v_i) : i \in \{1, \dots, k\}\}$ , and the edge weights are defined as  $x_{(v_{i-1}, v_i)} = a_i$  for  $i \in \{1, \dots, k\}$ . The node features are initialized as  $x_{v_0} = 0$  for the source node  $v_0$ , while all other nodes  $v_i$  for  $i > 0$  are initialized with  $x_{v_i} = \beta$ , representing an unreachable state.

## Appendix B. Warm-up: Single layer GNNs implement one step of BF

We start with the simple setting of a single layer GNN with shallow and narrow MLP components. This example provides key insights on why a perfectly (or almost perfectly) trained model can generalize. We analyze the general case of a multilayer GNN with wide and deep MLPs in Sec. C. Although the general case is more sophisticated technically, the approach follows similar intuitions. In particular, we later show that sparsity regularization can be used to reduce the analysis of GNN with wide and deep MLPs trained on a single BF step to the simple model analyzed in this section.

We start by proving Theorem 6, which shows how perfect accuracy on  $\mathcal{H}_{\text{small}}$  requires certain restrictions on parameters of the simple MinAgg GNN. Next, in Theorem 7 we show that such restrictions guarantee the parameters implement the BF algorithm. Finally, we extend this analysis to evaluate how MinAgg GNNs approximately minimizing the training loss perform on arbitrary graphs in Theorem 2

Suppose we have a *simple* 1-layer Bellman-Ford GNN,  $\mathcal{A}_\theta$ , where  $f^{\text{up},(0)} : \mathbb{R} \rightarrow \mathbb{R}$  and  $f^{\text{agg},(0)} : \mathbb{R}^2 \rightarrow \mathbb{R}$  are single layer MLPs. To be explicit:

$$h_u^{(1)} = \sigma(w_2 \min\{\sigma(W_1(x_v \oplus x_{(u,v)} + b_1)) : v \in \mathcal{N}(u) \cup \{u\}\} + b_2), \quad (12)$$

where  $\sigma$  is ReLU,  $W_1 \in \mathbb{R}^{1 \times 2}$ , and  $w_2, b_1, b_2 \in \mathbb{R}$ .

We consider the training set

$$\mathcal{H}_{\text{small}} = \{(P_1^{(0)}(a_i), P_1^{(1)}(a_i)) : i \in \{1, \dots, 4\}\} \cup \{(P_2^{(1)}(a_i, 0), P_2^{(2)}(a_i, 0)) : i \in \{5, \dots, 8\}\}. \quad (13)$$

For concreteness we take  $a_i = 2i$ , and we utilize these specific choices of edge weights in the proof of Theorem 2. However, any choice of  $a_i$  satisfying  $a_i \neq a_j$  if  $i \neq j$  and  $a_i > 0$  is sufficient for the other results in this section.

**Theorem 6** *If,  $\forall (H^{(t)}, \Gamma(H^{(t)})) \in \mathcal{H}_{\text{small}}$ ,*

$$\mathcal{A}_\theta(H^{(t)}) = \Gamma(H^{(t)}),$$

*i.e. the computed node features are  $h_u^{(1)}(H^{(t)}) = x_u(\Gamma(H^{(t)}))$  for all  $u \in V(H^{(t)})$ , then  $w_2 W_1 = \mathbf{1}$  and  $w_2 b_1 + b_2 = 0$ .*

**Proof** First, note that from the definition of  $P_1^{(0)}(a_i)$ , the source node is  $s = v_0$  so  $x_{v_0} = 0$ , and  $x_{v_1}(P_1^{(0)}(a_i)) = \beta$ . Additionally, given the training example  $(P_1^{(0)}(a_i), P_1^{(1)}(a_i)) \in \mathcal{H}_{\text{small}}$ , recall that  $P_1^{(0)}(a_i)$  is the input to  $\mathcal{A}_\theta$  (1-layer Bellman-Ford GNN). By the definition of  $\mathcal{A}_\theta$ , the computed node feature for  $v_1 \in V(P_1^{(0)}(a_i))$  is

$$\begin{aligned} h_{v_1}^{(1)} &= \sigma(w_2 \min\{\sigma(W_1(x_{v_1} \oplus x_{(v_1, v_1)})) + b_1), \sigma(W_1(x_{v_0} \oplus x_{(v_0, v_1)})) + b_1\}) + b_2 \\ &= \sigma(w_2 \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_i + b_1)\}) + b_2 \end{aligned}$$

where  $\sigma$  is the ReLU activation function. Since

$$\begin{aligned} \mathcal{A}_\theta(P_1^{(0)}(a_1)) &= P_1^{(1)}(a_1) \\ &\vdots \\ \mathcal{A}_\theta(P_1^{(0)}(a_4)) &= P_1^{(1)}(a_4), \end{aligned}$$

for each  $v_1 \in V(P_1^{(1)}(a_i))$ ,  $x_{v_1}(P_1^{(1)}(a_i)) = a_i$  so  $h_{v_1}^{(1)}(P_1^{(0)}(a_i)) = a_i$ . Therefore,

$$\begin{aligned} a_1 &= \sigma(w_2 \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_1 + b_1)\}) + b_2 \\ &\vdots \\ a_4 &= \sigma(w_2 \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_4 + b_1)\}) + b_2. \end{aligned}$$

Suppose  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_i + b_1)\}$  and  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_j + b_1)\}$  for  $i \neq j$ . Then  $a_i = a_j$  when  $i \neq j$  which is a contradiction. Therefore, there can be at most one  $i$  for which  $a_i = \sigma(w_2\sigma(W_{11}\beta + b_1) + b_2)$ . WLOG, assume that

$$a_i = \sigma(w_2\sigma(W_{12}a_i + b_1) + b_2)$$

where  $i \in [3]$ . Since  $a_i > 0$  and  $\sigma$  is the ReLU function, we have that  $a_i = w_2\sigma(W_{12}a_i + b_1) + b_2$  for  $i \in [3]$ . Suppose  $W_{12}a_i + b_1 \leq 0$  and  $W_{12}a_j + b_1 \leq 0$  for  $i, j \in [3]$  where  $i \neq j$ . Then  $a_i = a_j = b_2$  which is a contradiction. WLOG, assume that  $W_{12}a_i + b_1 > 0$  for  $i \in [2]$ . Then, we get the following system of linear equations

$$\begin{aligned} a_1 &= w_2W_{12}a_1 + w_2b_1 + b_2 \\ a_2 &= w_2W_{12}a_2 + w_2b_1 + b_2. \end{aligned}$$

These linear equations are only satisfied when  $w_2W_{12} = 1$  and  $w_2b_1 + b_2 = 0$ .

Now, consider  $\{(P_2^{(1)}(a_i, 0), P_2^{(2)}(a_i, 0)) : a_i \in \mathbb{R}^+, i \in \{5, \dots, 8\}, a_i \neq a_j\}$ . From the definition of  $P_2^{(1)}(a_i, 0)$ , we know that  $s = v_0$ ,  $x_{v_0}(P_2^{(1)}(a_i, 0)) = 0$ ,  $x_{v_1}(P_2^{(1)}(a_i, 0)) = a_i$ ,  $x_{v_2}(P_2^{(1)}(a_i, 0)) = \beta$ ,  $x_{(v_0, v_1)} = a_i$ , and  $x_{(v_1, v_2)} = 0$ . Since  $\mathcal{A}_\theta(P_2^{(1)}(a_i, 0)) = P_2^{(2)}(a_i, 0)$ , the computed node feature for  $v_2 \in V(P_2^{(1)}(a_i, 0))$  is

$$\begin{aligned} h_{v_2}^{(1)} &= a_i = \sigma(w_2 \min\{\sigma(W_1(x_{v_2} \oplus x_{(v_2, v_2)})) + b_1), \sigma(W_1(x_{v_1} \oplus x_{(v_1, v_2)})) + b_1\}) + b_2 \\ &= \sigma(w_2 \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{11}a_i + b_1)\}) + b_2 \end{aligned}$$

for  $i \in \{5, \dots, 8\}$ . Similar to above, we have that  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{11}a_i + b_1)\}$  can only occur for one  $i \in \{5, \dots, 8\}$ . Again, WLOG we can assume that  $a_8 = \sigma(w_2\sigma(W_{11}\beta + b_1) + b_2)$  and  $a_i = \sigma(w_2\sigma(W_{11}a_i + b_1) + b_2)$  for  $i \in \{5, 6, 7\}$ . Then, using a similar system of linear equations as above, we get that  $w_2W_{11} = 1$ .  $\blacksquare$

**Corollary 7** *Let  $\mathcal{A}_\theta$  be a simple 1-layer Bellman-Ford GNN, as given in Eq. (12). If  $\mathcal{A}_\theta(H^{(t)}) = \Gamma(H^{(t)})$  for all  $(H^{(t)}, \Gamma(H^{(t)})) \in \mathcal{H}_{\text{small}}$ , then for any  $G \in \mathcal{G}$  the MinAgg GNN outputs  $\mathcal{A}_\theta(G) = \Gamma(G)$  which means for any  $v \in V(G)$*

$$h_v^{(1)} = \min\{x_u + x_{(u,v)} : u \in \mathcal{N}(v)\}.$$

**Proof**

If  $\mathcal{A}_\theta(H^{(t)}) = \Gamma(H^{(t)})$  for all  $(H^{(t)}, \Gamma(H^{(t)})) \in \mathcal{H}_{\text{small}}$  then, by Theorem 6, we know that  $w_2W_{11} = \mathbb{1}$  and  $w_2b_1 + b_2 = 0$ . First, suppose  $w_2 < 0$ . Since  $w_2W_{11} = \mathbb{1}$ , we know that  $W_{11} = W_{12}$  and  $W_{11}, W_{12} < 0$ . Consider  $(P_1^{(0)}(a_i), P_1^{(1)}(a_i)) \in \mathcal{H}_{\text{small}}$ . Recall that  $a_i > 0$ . For any  $i \in \{1, \dots, 4\}$ , we have that  $v_1 \in V(P_1^{(0)}(a_i))$  gets the computed node feature

$$\begin{aligned} h_{v_1}^{(1)} &= \sigma(w_2 \min\{\sigma(W_{11}x_s + W_{12}x_{(s,v_1)} + b_1), \sigma(W_{11}x_{v_1} + W_{12}x_{(v_1,v_1)} + b_1)\} + b_2) \\ &= \sigma(w_2 \min\{\sigma(W_{11}a_i + b_1), \sigma(W_{11}\beta + b_1)\} + b_2) \end{aligned}$$

Since  $0 \leq a_i \ll \beta$ ,  $W_{11}\beta + b_1 \leq W_{11}a_i + b_1$  so

$$\min\{\sigma(W_{11}a_i + b_1), \sigma(W_{11}\beta + b_1)\} = \sigma(W_{11}\beta + b_1)$$

Then

$$h_{v_1}^{(1)} = \sigma(w_2\sigma(W_{11}\beta + b_1) + b_2)$$

for  $v_1 \in P_1^{(0)}(a_i)$  for any  $i \in \{1, \dots, 4\}$ . However, this is a contradiction because  $\mathcal{A}_\theta(P_1^{(0)}(a_i)) = P_1^{(1)}(a_i)$  for  $i \in \{1, \dots, 4\}$  and  $a_i \neq a_j$  for  $i \neq j$ . Therefore,  $w_2 > 0$  so  $W_{11}, W_{12} > 0$ .

Suppose  $w_2b_1 < 0$ . Because  $w_2 > 0$ ,  $b_1 < 0$ . Additionally, since  $w_2b_1 < 0$  and we know that  $w_2b_1 + b_2 = 0$ , we have that  $b_2 > 0$ . Then consider  $(P_1^{(0)}(a_1), P_1^{(1)}(a_1)) = (P_1^{(0)}(0), P_1^{(1)}(0)) \in \mathcal{H}_{\text{small}}$ . Then,  $v_1 \in V(P_1^{(0)}(0))$  gets the updated node feature

$$h_{v_1}^{(1)} = \sigma(w_2 \min\{\sigma(b_1), \sigma(W_{11}\beta + b_1)\} + b_2) = b_2.$$

This is contradiction because  $\mathcal{A}_\theta(P_1^{(0)}(a_1)) = P_1^{(1)}(a_1)$  which means that the computed node feature  $h_{v_1}^{(1)}$  should be  $a_1 = 0$ .

Now, given  $G^{(m)} \in \mathcal{G}$ , then given  $v \in V(G^{(m)})$ , the updated node feature for  $v \in V(\mathcal{A}_\theta(G^{(m)}))$  is

$$\begin{aligned}
 h_v^{(1)} &= \sigma(w_2 \min\{\sigma(W_{11}x_u + W_{12}x_{(v,u)} + b_1) : u \in \mathcal{N}(v)\} + b_2) \\
 &= \sigma(\min\{w_2\sigma(W_{11}(x_u + x_{(v,u)})) + b_1 : u \in \mathcal{N}(v)\} + b_2) \\
 &= \sigma(\min\{\sigma(w_2W_{11}(x_u + x_{(v,u)})) + w_2b_1 : u \in \mathcal{N}(v)\} + b_2) \text{ since } w_2 > 0 \\
 &= \sigma(\min\{\sigma(w_2W_{11}(x_u + x_{(v,u)})) + w_2b_1 + b_2 : u \in \mathcal{N}(v)\}) \\
 &= \sigma(\min\{\sigma(x_u + x_{(v,u)} + w_2b_1) + b_2 : u \in \mathcal{N}(v)\}) \\
 &= \sigma(\min\{x_u + x_{(v,u)} + w_2b_1 + b_2 : u \in \mathcal{N}(v)\}) \text{ since } x_u + x_{(v,u)} + w_2b_1 \geq 0 \\
 &= \sigma(\min\{x_u + x_{(v,u)} : u \in \mathcal{N}(v)\}) \\
 &= \min\{x_u + x_{(v,u)} : u \in \mathcal{N}(v)\}.
 \end{aligned}$$

■

**Lemma 8** Consider two points  $(x_1, y_1), (x_2, y_2) \in \mathbb{R}^2$  such that  $|x_1| < D$  and  $|x_2 - x_1| > 2$  and an affine function  $f(x) = ax + b$ . Suppose  $|f(x_1) - y_1| < \epsilon$  and  $|f(x_2) - y_2| < \epsilon$ . If  $a_0 = \frac{y_2 - y_1}{x_2 - x_1}$  and  $b_0 = y_1 - a_0x_1$  are the slope and y-intercept of a line passing through  $(x_1, y_1)$  and  $(x_2, y_2)$  then  $|a_0 - a| < \epsilon$  and  $|b_0 - b| < 2(1 + D)\epsilon$ .

**Proof** First,  $a = \frac{f(x_2) - f(x_1)}{x_2 - x_1}$  implies

$$\begin{aligned}
 |a_0 - a| &= \left| \frac{y_2 - y_1}{x_2 - x_1} - \frac{f(x_2) - f(x_1)}{x_2 - x_1} \right| \\
 &= \frac{1}{|x_2 - x_1|} |(y_2 - f(x_2)) - (y_1 - f(x_1))| \\
 &\leq \frac{1}{|x_2 - x_1|} (|y_2 - f(x_2)| + |y_1 - f(x_1)|) \\
 &\leq \frac{2\epsilon}{|x_2 - x_1|} \\
 &\leq \epsilon.
 \end{aligned}$$

Now, since  $b = f(x_1) - ax_1$  we have

$$\begin{aligned}
 |b_0 - b| &= |y_1 - a_0x_1 - (f(x_1) - ax_1)| \\
 &= |(y_1 - f(x_1)) - x_1(a_0 - a)| \\
 &\leq |y_1 - f(x_1)| + |x_1||a_0 - a| \\
 &< (1 + D)\epsilon.
 \end{aligned}$$

■

We now restate Theorem 2 with additional details and provide a proof.

**Theorem 9** Let  $0 < \epsilon < 1$ . If  $\forall (H^{(t)}, \Gamma(H^{(t)})) \in \mathcal{H}_{\text{small}}$ , a MinAgg GNN  $\mathcal{A}_\theta$  that, for  $u \in V(G^{(t)})$ , computes a node feature satisfying  $|h_u^{(1)}(G^{(t)}) - x_u(\Gamma(G^{(t)}))| < \frac{\epsilon}{20}$ . Then

(i)  $\|w_2 W_1 - \mathbf{1}\|_1 < \epsilon$  and  $|w_2 b_1 + b_2| < 20\epsilon$

(ii)  $w_2, W_{11}, W_{12} \geq 0$

(iii) For  $G \in \mathcal{G}$  and  $v \in V(G)$

$$(1 - \epsilon)x_v(G) - \epsilon \leq h_v^{(1)}(G) \leq (1 + \epsilon)x_v(G) + \epsilon$$

**Proof**

(i) We first show part (i) i.e. if  $|h_u^{(1)}(G^{(t)}) - x_u(\Gamma(G^{(t)}))| < \frac{\epsilon}{20}$ , for any  $(G^{(t)}, \Gamma(G^{(t)})) \in \mathcal{H}_{\text{small}}$ , then  $\|w_2 W_1 - \mathbf{1}\|_1 < \epsilon$  and  $|w_2 b_1 + b_2| < 20\epsilon$ . Let  $\epsilon_0 = \frac{\epsilon}{20}$ . Given the definition of  $\mathcal{A}_\theta$ , the computed node feature for  $v_1 \in V(P_1^{(0)}(a_i))$  for  $i \in \{1, \dots, 4\}$  is

$$h_{v_1}^{(1)} = \sigma[w_2 \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_i + b_1)\}]$$

Since  $|h_{v_1}^{(1)}(P_1^{(0)}(a_i)) - x_{v_1}(P_1^{(1)}(a_i))| < \epsilon_0$ ,

$$|\sigma(w_2 \min(\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_1 + b_1)) + b_2) - a_1| < \epsilon_0$$

$\vdots$

$$|\sigma(w_2 \min(\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_4 + b_1)) + b_2) - a_4| < \epsilon_0$$

Suppose  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_i + b_1)\}$  and  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_j + b_1)\}$  for  $i \neq j$ . Then  $|\sigma(w_2\sigma(W_{11}\beta + b_1) + b_2) - a_i| < \epsilon_0$  and  $|\sigma(w_2\sigma(W_{11}\beta + b_1) + b_2) - a_j| < \epsilon_0$  so  $|a_i - a_j| < 2\epsilon_0 < 2$ . This is a contradiction because for any  $i \neq j$ ,  $|a_i - a_j| \geq 2$ .

Suppose that  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_{i_1} + b_1)\}$  and  $\sigma(W_{11}\beta + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_{i_2} + b_1)\}$  for  $i_1, i_2 \in \{1, 2, 3, 4\}$  for  $i_1 \neq i_2$ . This implies that  $|a_{i_1} - a_{i_2}| < 2$  which is a contradiction. Thus, w.l.o.g. we can assume that for  $i \in \{1, 2, 3\}$ ,  $\sigma(W_{12}a_i + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_i + b_1)\}$ . Additionally, suppose  $W_{12}a_i + b_1 < 0$  and  $W_{12}a_j + b_1 < 0$  for  $i \neq j$  and  $i, j \in \{1, 2, 3\}$ . Then,  $h_{v_1}(P_1^{(0)}(a_i)) = \sigma(b_2)$  and  $h_{v_1}(P_1^{(0)}(a_j)) = \sigma(b_2)$  so  $|\sigma(b_2) - a_i| < \epsilon_0$  and  $|\sigma(b_2) - a_j| < \epsilon_0$ . From here, we get that  $|a_i - a_j| < 2$ , which is a contradiction. Therefore, we assume that  $\sigma(W_{12}a_i + b_1) = \min\{\sigma(W_{11}\beta + b_1), \sigma(W_{12}a_{i_1} + b_1)\}$  and  $W_{12}a_i + b_1 > 0$  for  $i = \{1, 2\}$ .

Then we have

$$|(a_1 w_2 W_{12} + w_2 b_1 + b_2) - a_1| \leq \epsilon_0$$

and

$$|(a_2 w_2 W_{12} + w_2 b_1 + b_2) - a_2| \leq \epsilon_0.$$

Note that  $f(a) = a w_2 W_{12} + w_2 b_1 + b_2$  is an affine function with slope  $m = w_2 W_{12}$  and intercept  $w_2 b_1 + b_2$ . As  $|a_1|, \dots, |a_4| < 9$  and  $|a_i - a_j| > 2$ , by Lemma 8,

$$|w_2 W_{11} - 1| < \epsilon_0 < \frac{\epsilon}{2}$$

$$|w_2 b_1 + b_2| < 2 \cdot (1 + 9)\epsilon_0 = 20\epsilon_0 = \epsilon.$$

A parallel method using  $\mathcal{H}_{\text{small}}$  with  $a_i = 2i$  for  $i \in \{5, 6, 7, 8\}$  follows for bounding  $|w_2 W_{11} - 1| < \frac{\epsilon}{2}$ .

- (ii) Now, we turn our attention to part (ii) and show that  $w_2, W_{11}, W_{12} \geq 0$ . Suppose  $w_2 < 0$ , which implies  $W_{11}, W_{12} < 0$  as otherwise  $w_2 W_{11}, w_2 W_{12} < 0$  and  $|w_2 W_1 - 1| > \epsilon$  (recall  $0 < \epsilon < 1$ ). The computed node feature for  $v_1 \in V(P_1^{(0)}(a_i))$  is then

$$\begin{aligned} h_{v_1}^{(1)} &= \sigma(w_2 \min\{\sigma(W_{11}x_s + W_{12}x_{(s,v_1)} + b_1), \sigma(W_{11}x_{v_1} + W_{12}x_{(v_1,v_1)} + b_1)\} + b_2) \\ &= \sigma(w_2 \min\{\sigma(W_{11}a_i + b_1), \sigma(W_{11}\beta + b_1)\} + b_2) \\ &= \sigma(w_2 \sigma(W_{11}\beta + b_1)) + b_2. \end{aligned}$$

Note that the above inequality follows from the fact that  $W_{11}a_i + b_1 > W_{11}\beta + b_1$  since  $a_i < \beta$  and  $W_{11} < 0$ . For  $v_1 \in V(P_1^{(0)}(a_i))$  for all  $i \in \{1, 2, 3, 4\}$ ,  $h_{v_1}^{(1)} = \sigma(w_2 \sigma(W_{11}\beta + b_1)) + b_2$ . However, this is a contradiction, since  $|a_i - a_j| \geq 2$  for all  $i \neq j$  where  $i, j \in \{1, \dots, 4\}$ . Thus,  $w_2, W_{11}, W_{12} \geq 0$ .

- (iii) We will now show that given any  $G \in \mathcal{G}$  and  $v \in V(G)$ , a neural network with the weights given in part (i) will approximately yield a single step of Bellman-Ford i.e.

$$(1 - \epsilon)h_v^{(1)}(G) - \epsilon \leq h_v^{(1)}(G) \leq (1 + \epsilon)x_v(G) + \epsilon.$$

Since  $\|w_2 W_1 - 1\|_1 < \epsilon$  from part (i), we know that  $|w_2 W_{11} - 1| < \epsilon$  and  $|w_2 W_{12} - 1| < \epsilon$ . Additionally, we know that for any  $G \in \mathcal{G}$  and  $v \in V(G)$ ,

$$h_v^{(1)}(G) = \sigma(w_2 \min\{\sigma(W_{11}x_u + W_{12}x_{(u,v)} + b_1) : u \in \mathcal{N}(v)\} + b_2)$$

From (ii), we know that  $W_{11} \geq 0$ ,  $W_{12} \geq 0$ , and  $W_{11}x_u + W_{12}x_{(u,v)} + b_1 > 0$  so the ReLU activation function  $\sigma$  can be removed from the aggregation MLP i.e.,

$$h_v^{(1)}(G) = \sigma(w_2 \min\{W_{11}x_u + W_{12}x_{(u,v)} + b_1 : u \in \mathcal{N}(v)\} + b_2).$$

Suppose

$$u' = \operatorname{argmin}_{u \in \mathcal{N}(v)} \{W_{11}x_u + W_{12}x_{(u,v)} + b_1\}$$

and

$$u^* = \operatorname{argmin}_{u \in \mathcal{N}(v)} \{x_u + x_{(u,v)}\}.$$

Note that  $x_v(\Gamma(G)) = x_{u^*} + x_{(u^*,v)}$ . Then,

$$\begin{aligned} h_v^{(1)}(G) &= \sigma(w_2(W_{11}x_{u'} + W_{12}x_{(u',v)} + b_1) + b_2) \\ &\leq \sigma(w_2 W_{11}x_{u^*} + w_2 W_{12}x_{(u^*,v)} + w_2 b_1 + b_2) \\ &\leq \sigma((1 + \epsilon)(x_{u^*} + x_{(u^*,v)}) + w_2 b_1 + b_2) \end{aligned}$$

Note that if  $w_2 b_1 + b_2 \leq 0$ , then

$$h_v^{(1)}(G) \leq \sigma((1 + \epsilon)(x_{u^*} + x_{(u^*,v)}) + w_2 b_1 + b_2) \leq \sigma((1 + \epsilon)(x_{u^*} + x_{(u^*,v)})) = (1 + \epsilon)(x_{u^*} + x_{(u^*,v)}).$$

If  $w_2b_1 + b_2 > 0$ , then

$$\begin{aligned} h_v^{(1)}(G) &\leq \sigma((1 + \epsilon)(x_{u^*} + x_{(u^*,v)}) + w_2b_1 + b_2) \\ &\leq (1 + \epsilon)(x_{u^*} + x_{(u^*,v)}) + \epsilon \\ &= (1 + \epsilon)x_v(\Gamma(G)) + \epsilon. \end{aligned}$$

In both cases,  $h_v^{(1)}(G) \leq (1 + \epsilon)x_v(\Gamma(G)) + \epsilon$ .

Now, we consider the lower bound and show that  $(1 - \epsilon)x_v(\Gamma(G)) - \epsilon < h_v^{(1)}(G)$ . By the definition of  $u^*$ ,

$$x_{u^*} + x_{u^*,v} \leq x_{u'} + x_{(u',v)}$$

Note that because  $0 < \epsilon < 1$ , we have  $0 < 1 - \epsilon < 1$  and  $\frac{1}{1-\epsilon} > 1$ . We will consider two cases: when  $w_2b_1 + b_2 \geq 0$  and when  $w_2b_1 + b_2 < 0$ . Let  $w_2b_1 + b_2 > 0$ . Then

$$\begin{aligned} x_{u^*} + x_{u^*,v} &\leq x_{u'} + x_{(u',v)} + w_2b_1 + b_2 \\ &\leq \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{u'} + \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{(u',v)} + \left(\frac{1-\epsilon}{1-\epsilon}\right)(w_2b_1 + b_2) \\ &\leq \left(\frac{1}{1-\epsilon}\right) \cdot (1-\epsilon)x_{u'} + \left(\frac{1}{1-\epsilon}\right) \cdot (1-\epsilon)x_{(u',v)} + \left(\frac{1}{1-\epsilon}\right)(w_2b_1 + b_2) \\ &\leq \left(\frac{1}{1-\epsilon}\right) \left( (1-\epsilon)x_{u'} + (1-\epsilon)x_{(u',v)} + w_2b_1 + b_2 \right) \\ &\leq \left(\frac{1}{1-\epsilon}\right) \cdot (w_2W_{11}x_{u'} + w_2W_{12}x_{(u',v)} + w_2b_1 + b_2) \\ &= \left(\frac{1}{1-\epsilon}\right) \cdot h_v(G) \end{aligned}$$

Therefore,

$$(1 - \epsilon)(x_{u^*} + x_{u^*,v}) = (1 - \epsilon)x_v(\Gamma(G)) \leq h_v(G).$$

Let  $w_1b_1 + b_2 < 0$ . We know that

$$w_1W_{11}x_{u'} + w_2W_{12}x_{(u',v)} + w_2b_1 + b_2 \leq w_1W_{11}x_{u'} + w_1W_{12}x_{(u',v)} + w_2b_1 + b_2 + \epsilon$$

Since  $|w_1b_1 + b_2| < \epsilon$ ,  $w_1b_1 + b_2 + \epsilon > 0$ . Therefore,

$$\begin{aligned} x_{u^*} + x_{u^*,v} &\leq x_{u'} + x_{(u',v)} + w_2b_1 + b_2 + \epsilon \\ &\leq \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{u'} + \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{(u',v)} + \left(\frac{1-\epsilon}{1-\epsilon}\right)((w_2b_1 + b_2) + \epsilon) \\ &\leq \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{u'} + \left(\frac{1-\epsilon}{1-\epsilon}\right)x_{(u',v)} + \left(\frac{1}{1-\epsilon}\right)((w_2b_1 + b_2) + \epsilon) \\ &\leq \left(\frac{1}{1-\epsilon}\right)(h_v(G) + \epsilon) \end{aligned}$$

Thus,

$$(1 - \epsilon)(x_{u^*} + x_{u^*,v}) - \epsilon \leq h_v(G)$$

■

### Appendix C. Sparsity regularized deep GNNs implement BF

In this section we analyze GNNs that are large both in their number of layers and the size of their respective MLPs. The key to showing these complex GNNs implement the BF algorithm is the introduction of sparsity regularization to the loss. With this type of regularization we can show any solution that approximates the global minimum must have only a few non-zero parameters. Furthermore, any GNN with so few non-zero parameters can solve shortest path problems only via the BF algorithm. In short, although the model is over-parameterized, solutions approximating the global minimum are not.

Our overarching approach is as follows. We first give an implementation of BF by GNN with a small number of non-zero parameters  $S$ . Next, we show that, on our constructed training set, any GNN with less than  $S$  non-zero parameters has large error. This allows us to conclude that the global minimum of the sparsity regularized loss must have exactly  $S$  non-zero parameters. This sparsity allows us to simplify the MinAgg GNN update to include only a few parameters. We then derive approximations to these parameters which show the MinAgg GNN must be implementing BF algorithm, up to some scaling factor.

A key strategy in this section is to track the dependencies of the functions  $f^{\text{agg},(\ell)}$  and  $f^{\text{up},(\ell)}$  on their components. In particular, we say a function  $f$  depends on a component or set of components if it is not constant over these components. Note that inputs to these functions are always non-negative (they are always preceded by a ReLU), so by constant we mean constant over all non-negative values. The precise definitions are as follows.

**Definition 10** For  $\ell \in [L]$  the function  $f^{\text{agg},(\ell)}$  **depends on its node component** iff it is not constant over its first  $d_{\ell-1}$  components, i.e., there exists  $x, y \in \mathbb{R}_{\geq 0}^{d_{\ell-1}+1}$  with  $x \neq y$  and  $x_{d_{\ell-1}+1} = y_{d_{\ell-1}+1}$  such that

$$f^{\text{agg},(\ell)}(x) \neq f^{\text{agg},(\ell)}(y).$$

The function  $f^{\text{agg},(\ell)}$  **depends on its edge component** iff it is not constant over its edge component (the  $(d_{\ell-1} + 1)$ th component). That is, there exists  $x, y \in \mathbb{R}_{\geq 0}^{d_{\ell-1}+1}$  with  $x \neq y$  and  $x_i = y_i$  for  $i \in \{1, \dots, d_{\ell-1}\}$  such that

$$f^{\text{agg},(\ell)}(x) \neq f^{\text{agg},(\ell)}(y).$$

**Definition 11** For  $\ell \in [L]$  the function  $f^{\text{up},(\ell)}$  **depends on its aggregation component** iff it is not constant over its first  $d$  components, i.e., there exists  $x, y \in \mathbb{R}_{\geq 0}^{d+d_{\ell-1}}$  with  $x \neq y$  and  $x_i = y_i$  for  $i \in \{d+1, \dots, d_{\ell-1}\}$  such that

$$f^{\text{up},(\ell)}(x) \neq f^{\text{up},(\ell)}(y).$$

The function  $f^{\text{up},(\ell)}$  **depends on its skip component** iff it is not constant over its last  $d_{\ell-1}$  components, i.e., there exists  $x, y \in \mathbb{R}_{\geq 0}^{d+d_{\ell-1}}$  with  $x \neq y$  and  $x_i = y_i$  for  $i \in \{1, \dots, d\}$  such that

$$f^{\text{up},(\ell)}(x) \neq f^{\text{up},(\ell)}(y).$$

Our approach proceeds by showing that requisite dependencies can only be achieved if there is some minimal number of non-zero entries in  $\theta$ .

### C.1. Implementing BF

We begin by showing there is a choice of parameters that makes the MinAgg GNN implement  $K$  steps of the BF algorithm.

**Lemma 12** *Let  $L \geq K > 0$ . For an  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$  with  $m$ -layer update and aggregation MLPs and parameters  $\theta$ , there is an assignment of  $\theta$  with  $mL + mK + K$  non-zero values such that  $\mathcal{A}_\theta$  implements  $K$  steps of the BF algorithm, i.e., for any  $G \in \mathcal{G}$*

$$\mathcal{A}_\theta(G) = \Gamma^K(G).$$

**Proof** We proceed by assigning parameters to  $\mathcal{A}_\theta$  such that  $\mathcal{A}_\theta$  simulates  $K$  steps of Bellman-Ford, i.e., for any  $G \in \mathcal{G}$ ,  $\mathcal{A}_\theta(G) = \Gamma^K(G)$ . For  $\ell \in [L]$ , let  $f^{\text{agg},(\ell)} : \mathbb{R}^{d_{\ell-1}+1} \rightarrow \mathbb{R}^d$  and  $f^{\text{up},(\ell)} : \mathbb{R}^{d+d_{\ell-1}} \rightarrow \mathbb{R}^d$  be the  $m$ -layer update and aggregation MLPs respectively. Note that  $d_0 = d_L = 1$ , and for  $\ell \in \{1, \dots, L-1\}$  the hidden layer dimension is  $d_\ell = d$ , for some arbitrary  $d \geq 1$ . The parameters for  $f^{\text{agg},(\ell)}$  and  $f^{\text{up},(\ell)}$  are  $\{(W_j^{\text{agg},(\ell)}, b_j^{\text{agg},(\ell)})\}_{j \in [m]}$  and  $\{(W_j^{\text{up},(\ell)}, b_j^{\text{up},(\ell)})\}_{j \in [m]}$ , respectively, where for  $j \in [m]$ ,

$$\begin{aligned} W_j^{\text{agg},(\ell)} &\in \mathbb{R}^{d_j^{\text{agg},(\ell)} \times d_{j-1}^{\text{agg},(\ell)}} \\ b_j^{\text{agg},(\ell)} &\in \mathbb{R}^{d_j^{\text{agg},(\ell)}} \\ W_j^{\text{up},(\ell)} &\in \mathbb{R}^{d_j^{\text{up},(\ell)} \times d_{j-1}^{\text{up},(\ell)}} \\ b_j^{\text{up},(\ell)} &\in \mathbb{R}^{d_j^{\text{up},(\ell)}}. \end{aligned}$$

The dimension of these parameters are

$$\begin{aligned} d_j^{\text{agg},(\ell)} &= \begin{cases} d_{\ell-1} + 1 & \text{if } j = 0 \\ d & \text{otherwise} \end{cases} \\ d_j^{\text{up},(\ell)} &= \begin{cases} d_{\ell-1} + d & \text{if } j = 0 \\ d_\ell & \text{if } j = m \\ d & \text{otherwise} \end{cases}. \end{aligned}$$

Now we give values of these parameters that make  $\mathcal{A}_\theta$  implement the BF algorithm. Let  $b_j^{\text{up},(\ell)} = \mathbf{0}$  and  $b_j^{\text{agg},(\ell)} = \mathbf{0}$  for all  $\ell \in [L]$  and  $j \in [m]$ . Set

$$W_1^{\text{agg},(1)} = \begin{pmatrix} 1 & 1 \\ \vdots & \vdots \\ 0 & 0 \end{pmatrix}$$

and for  $\ell \in \{2, \dots, K\}$

$$W_1^{\text{agg},(\ell)} = \begin{pmatrix} 1 & 0 & \dots & 0 & 1 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & & & \vdots \\ 0 & 0 & \dots & & 0 \end{pmatrix}. \quad (14)$$

This choice makes  $W_1^{\text{agg},(\ell)}$  sum the edge weight and first component of the node feature into the first component of the resulting vector. That is,  $[W_1^{\text{agg},(\ell)}(h_v^{(\ell-1)} \oplus x_{(v,u)})]_1 = [h_v^{(\ell-1)}]_1 + x_{(v,u)}$ . Next, set

$$W_j^{\text{agg},(\ell)} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & & \\ 0 & \dots & & 0 \end{pmatrix} \quad \text{for } \ell \in [K] \text{ and } j \in \{2, \dots, m\},$$

$$W_j^{\text{up},(\ell)} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & & \\ 0 & \dots & & 0 \end{pmatrix} \quad \text{for } \ell \in [K] \text{ and } j \in [m].$$

Finally, for  $\ell \in \{K+1, \dots, L\}$ , let

$$W_j^{\text{agg},(\ell)} = \mathbf{0} \quad \text{for } j \in [m]$$

$$W_1^{\text{up},(\ell)} = \begin{pmatrix} 0 & 0 & \dots & 1 \\ 0 & \ddots & & \vdots \\ \vdots & & & \\ 0 & \dots & & 0 \end{pmatrix}$$

$$W_j^{\text{up},(\ell)} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & & \\ 0 & \dots & & 0 \end{pmatrix} \quad \text{for } j \in \{2, \dots, m\}.$$

Given the above assignments of  $W_j^{\text{agg},(\ell)}$  and  $W_j^{\text{up},(\ell)}$ , the final  $L - K$  layers implement the identity on the first component of the node feature, i.e.,  $[h_v^{(\ell)}]_1 = [h_v^{(\ell-1)}]_1$  for  $\ell \in \{K+1, \dots, L\}$ . Since edge weights are always non-negative and there are no negative parameters in the above, the ReLU activations can be ignored. Then, for all  $v \in V$  and for  $\ell \leq K$ , we get

$$[h_v^{(\ell)}]_1 = \min\{[h_u^{(\ell-1)}]_1 + x_{(u,v)} \mid u \in \mathcal{N}(v)\}$$

which is the BF algorithm update. This implies, by the correctness of the BF algorithm, that  $[h_v^{(K)}]_1$  is the  $K$ -step shortest path distance. The last  $L - K$  layers of the GNN implement the identity function so  $h_v^{(L)} = [h_v^{(L)}]_1 = [h_v^{(K)}]_1$  is also the  $K$ -step shortest path distances. Perfect accuracy is then achieved on all  $K$ -step shortest path instances. ■

We later show that the requirement  $L \geq K$  is indeed necessary (Corollary 22).

## C.2. Training set

Our training set is comprised of multiple parts, which we describe in this subsection. The first set of training instances is used to regulate how the MinAgg GNN scales features throughout computation.

**Definition 13 ()** For  $k \in [K]$ , define  $\mathcal{H}_{k,K}$  as

$$\mathcal{H}_{k,K} = \{P_{K+1}^{(1)}(a, 0, \dots, 0, b, 0, \dots, 0) : (a, b) \in \{0, 1, \dots, 2K\} \times \{0, 2K + 1\}\}$$

where  $P_{K+1}^{(1)}(a, 0, \dots, 0, b, 0, \dots, 0)$  is the attributed  $K$ -edge path graph with weight  $a$  for the first edge, weight  $b$  for the  $(k + 1)$ th edge, and weight zero for all other edges.

Next, we define graph that is used to show the MinAgg GNN must have at least  $K$  steps that depend on both edge weights and neighboring node features. If these conditions are not met, then the MinAgg GNN is not expressive enough to compute the shortest path distances in this graph.

**Definition 14 ()** Let  $H^{(0),K}$  be a 0-step BF instance with  $2K + 2$  vertices

$$V = \{v_0, v_1, \dots, v_K\} \cup \{u_0, u_1, \dots, u_K\},$$

edges

$$E = \{(v_{i-1}, v_i) \mid i \in [K]\} \cup \{(u_{i-1}, u_i) \mid i \in [K]\} \cup \{(u_{i-1}, v_i) \mid i \in [K]\} \cup \{(v_{i-1}, u_i) \mid i \in [K]\},$$

edge features  $X_e$  given by

$$x_{(w,q)} = \begin{cases} 1 & \text{if } (w, q) = (u_{k-1}, v_k) \text{ or } (w, q) = (v_{k-1}, u_k) \text{ for } k \in [K] \\ 0 & \text{otherwise} \end{cases},$$

and initial node features  $X_v$  given by

$$x_w = \begin{cases} 0 & \text{if } w = v_0 \\ \beta & \text{otherwise} \end{cases}.$$

We also write  $H_K^{(K)} = \Gamma^K(H_K^{(0)})$ .

The complete training set also includes  $P_1^{(0)}(1), P_2^{(1)}(1, 0)$ .

**Definition 15** For  $K > 1$ , we let

$$\begin{aligned} \mathcal{G}_{scale,K} &= \cup_{K \geq k > 1} \mathcal{H}_{k,K} \\ \mathcal{G}_K &= \mathcal{G}_{scale,K} \cup \{P_1^{(0)}(1), P_2^{(1)}(1, 0), H_K^{(0)}\} \\ \mathcal{G}_K &= \{(G^{(t)}, \Gamma(G^{(t)})) : G^{(t)} \in \mathcal{G}_K\}. \end{aligned}$$

Note the distinction here between  $\mathcal{G}_K$  which is a set of graphs and  $\mathcal{G}_K$ , which contains pairs of graphs (an input graph and a target graph).

### C.3. Sparsity structure

Here we show that the training set  $\mathcal{G}_K$  requires a MinAgg GNN to have a minimal sparsity to achieve good performance. Furthermore, this non-zero parameters must follow a particular structure.

**Definition 16** An *isomorphism* between two attributed graphs  $G = (V, E, X_v, X_e)$  and  $G' = (V', E', X'_v, X'_e)$  is a bijection  $\phi : V \rightarrow V'$  satisfying

$$(u, v) \in E \text{ if and only if } (\phi(u), \phi(v)) \in E'$$

and

$$\begin{aligned} x_v &= x'_{\phi(v)} \quad \forall v \in V \\ x_{(v,u)} &= x'_{(\phi(v), \phi(u))} \quad \forall (v, u) \in E. \end{aligned}$$

**Fact 17** Suppose  $\phi$  is an isomorphism between two attributed graphs  $G = (V, E, X_v, X_e)$  and  $G' = (V', E', X'_v, X'_e)$ . Then  $\phi$  is also an isomorphism between  $A_\theta(G)$  and  $A_\theta(G')$ .

**Definition 18** For an  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$ , we say that a layer  $\ell \in [L]$  is **message passing** if the aggregation function  $f^{\text{agg},(\ell)}$  depends on its node component. A **edge-dependent message passing layer** is a message passing layer for which  $f^{\text{agg},(\ell)}$  also depends on its edge component. A layer is **stationary** if it is not message passing

The update at each node can only depend on neighboring node features through the node component of the aggregation function. Thus, for stationary layers, each updated node feature only depends on the previous value of the node feature and the value of adjacent edges.

**Fact 19** Consider a MinAgg GNN  $\mathcal{A}_\theta$  such that its  $\ell$ th layer is stationary. Then, taking  $\theta$  to be fixed, the feature  $h_v^{(\ell)}$  is only a function of  $h_v^{(\ell-1)}$  and  $x_{(u,v)}$  for  $u \in \mathcal{N}(v)$ .

The importance of message passing layers is that these are the only layers for which an updated node depends on the previous features of its neighboring nodes. This is made precise by the following statement.

**Claim 20** Consider a MinAgg GNN  $\mathcal{A}_\theta$  acting on a graph  $G = (V, E, X_v, X_e)$  and consider two nodes  $v, w \in V$  such that  $v$  is  $j$  steps away from  $w$ . Suppose for some  $\ell \in [L]$  there are  $k$  layers in  $[\ell]$  that are message passing with  $k < j$ . Then the feature  $h_v^{(\ell)}(G)$  is independent of  $x_w(G)$ . That is, for any graph  $G' = (V, E, X_v, X'_e)$  that differs from  $G$  only in the feature  $x_w(G')$ , we have  $h_v^{(\ell)}(G) = h_v^{(\ell)}(G')$ . Similarly, for any edge  $(u, w) \in E$  if both  $u$  and  $w$  are  $j$  steps away from  $v$  then  $h_v^{(\ell)}(G)$  is independent of  $x_{(u,w)}(G)$ .

**Proof** We proceed by induction so assume the statement holds for  $j - 1$ . Note the base case of  $j = 1$  is immediate from Fact 19 as this implies no message passing has occurred. Suppose  $v$  is  $j$  steps away from  $w$ . Let  $\ell'$  be the largest  $\ell$  in  $\{1, \dots, \ell\}$  that is message passing. By definition of the MinAgg GNN,

$$h_v^{(\ell')} = f^{\text{up},(\ell')} \left( \min \{ f^{\text{agg},(\ell')} (h_u^{(\ell'-1)} \oplus x_{(u,v)}) : u \in \mathcal{N}(v) \} \right). \quad (15)$$

Every node in  $\mathcal{N}(v)$  is at least  $j - 1$  steps from  $w$  and there are at most  $k - 1$  message passing steps in  $[\ell' - 1]$ . Invoking the inductive hypothesis for  $j - 1$  yields that  $h_u^{(\ell'-1)}$  for  $u \in \mathcal{N}(v)$  is

independent of  $x_u(G)$ . Thus, since every variable in the expression for  $h_v^{(\ell')}$  is independent of  $x_u(G)$ , the feature  $h_v^{(\ell')}$  is independent as well. Finally, if  $h_v^{(\ell')}$  is independent of  $x_u(G)$  then by Fact 19,  $h_v^{(\ell)}$  is also independent of  $x_u(G)$  since all the layers between  $\ell'$  and  $\ell$  are stationary. (If  $\ell$  is message passing then  $\ell' = \ell$ .)

A parallel argument can be used to show independence of  $h_v^{(\ell)}$  from  $x_{(u,w)}$  in the case that  $v$  is  $j$  steps away from both  $u$  and  $w$  ■

**Lemma 21** *An  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$  satisfies*

$$h_{v_K}^{(L)}(H_K^{(0)}) \neq h_{u_K}^{(L)}(H_K^{(0)})$$

*only if it has at least  $K$  edge-dependent message passing layers.*

**Proof**

We proceed by proving claim (\*), regarding the output of the MinAgg GNN on  $H_K^{(0)}$ . This claim is shown through two cases. Note that since in this proof only the graph  $H_K^{(0)}$  is considered, we suppress notation referring to the graph for simplicity.

Claim (\*) Let  $k' \in [K]$  and  $\ell \in [L]$ . If for all  $k \geq k'$ ,

$$h_{v_k}^{(\ell-1)} = h_{u_k}^{(\ell-1)},$$

and  $f^{\text{agg},(\ell)}$  is not edge-dependent message passing, then for all  $k \geq k'$ ,

$$h_{v_k}^{(\ell)} = h_{u_k}^{(\ell)}.$$

If  $f^{\text{agg},(\ell)}$  is not edge-dependent message passing, then it either does not depend on its node component or does not depend on its edge component.

Case I:  $f^{\text{agg},(\ell)}$  **does not depend on its node component.** For  $k \geq k'$ ,

$$\begin{aligned} h_{v_k}^{(\ell)} &= f^{\text{up},(\ell)} \left( \min\{f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus x_{(w,v_k)}) : w \in \{v_{k-1}, u_{k-1}, v_k, v_{k+1}, u_{k+1}\}\} \oplus h_{v_k}^{(\ell-1)} \right) \\ &= f^{\text{up},(\ell)} \left( \min\{f^{\text{agg},(\ell)}(\cdot \oplus 0), f^{\text{agg},(\ell)}(\cdot \oplus 1)\} \oplus h_{v_k}^{(\ell-1)} \right) \end{aligned}$$

where a center dot  $\cdot$  is used to indicate the lack of dependence on the node feature. (Any value can replace  $\cdot$  and the expression does not change since  $f^{\text{agg},(\ell)}$  is constant over its node component.) The key here is that since  $f^{\text{agg},(\ell)}$  does not depend on its node component, only the set of edge weights incident on  $v_k$ , which is  $\{0, 1\}$ , effects the value of

$$\min\{f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus x_{(w,v_k)}) : w \in \{v_{k-1}, u_{k-1}, v_k, v_{k+1}, u_{k+1}\}\}.$$

Similarly, for  $u_k$  we have

$$\begin{aligned} h_{u_k}^{(\ell)} &= f^{\text{up},(\ell)} \left( \min\{f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus x_{(w,u_k)}) : w \in \{v_{k-1}, u_{k-1}, u_k, v_{k+1}, u_{k+1}\}\} \oplus h_{u_k}^{(\ell-1)} \right) \\ &= f^{\text{up},(\ell)} \left( \min\{f^{\text{agg},(\ell)}(\cdot \oplus 0), f^{\text{agg},(\ell)}(\cdot \oplus 1)\} \oplus h_{u_k}^{(\ell-1)} \right) \end{aligned}$$

and since  $h_{v_k}^{(\ell-1)} = h_{v_k}^{(\ell-1)}$ , we get  $h_{u_k}^{(\ell)} = h_{v_k}^{(\ell)}$ .

Case II:  $f^{\text{agg},(\ell)}$  **does not depend on its edge component.**

For  $k \geq k'$ ,

$$\begin{aligned} h_{v_k}^{(\ell)} &= f^{\text{up},(\ell)} \left( \min \{ f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus x_{(w,v_k)}) : w \in \{v_{k-1}, u_{k-1}, v_k, v_{k+1}, u_{k+1}\} \} \oplus h_{v_k}^{(\ell-1)} \right) \\ &= f^{\text{up},(\ell)} \left( \min \{ f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus \cdot) : w \in \{v_{k-1}, u_{k-1}, v_k, v_{k+1}, u_{k+1}\} \} \oplus h_{v_k}^{(\ell-1)} \right) \end{aligned}$$

and

$$\begin{aligned} h_{u_k}^{(\ell)} &= f^{\text{up},(\ell)} \left( \min \{ f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus x_{(w,u_k)}) : w \in \{v_{k-1}, u_{k-1}, u_k, v_{k+1}, u_{k+1}\} \} \oplus h_{u_k}^{(\ell-1)} \right) \\ &= f^{\text{up},(\ell)} \left( \min \{ f^{\text{agg},(\ell)}(h_w^{(\ell-1)} \oplus \cdot) : w \in \{v_{k-1}, u_{k-1}, u_k, v_{k+1}, u_{k+1}\} \} \oplus h_{u_k}^{(\ell-1)} \right). \end{aligned}$$

However, since  $h_{v_k}^{(\ell-1)} = h_{u_k}^{(\ell-1)}$ , the minimums in the expressions for  $h_{v_k}^{(\ell)}$  and  $h_{u_k}^{(\ell)}$  are taken over the same set and so  $h_{v_k}^{(\ell)} = h_{u_k}^{(\ell)}$ .

Now that we have proved claim (\*) we proceed by induction to show that for any  $\ell' \in [L]$  if there are  $k'$  edge-dependent message passing layers in  $[\ell']$  then  $h_{u_k}^{(\ell')} = h_{v_k}^{(\ell')}$  for all  $k > k'$ . This holds for the base case of  $k' = 0$  since at initialization  $h_{v_k}^{(0)} = h_{u_k}^{(0)}$  for all  $k \in [K]$ , and (\*) yields that  $h_{v_k}^{(\ell')} = h_{u_k}^{(\ell')}$  since no layer in  $[\ell']$  is edge-dependent message passing. Now suppose the inductive hypothesis holds for  $k' - 1$  edge-dependent message passing layers. That is, suppose that for any  $\ell' \in [L]$  if there are  $k' - 1$  edge-dependent message passing layers in  $[\ell']$  then  $h_{u_k}^{(\ell')} = h_{v_k}^{(\ell')}$  for all  $k > k' - 1$ . Now suppose for some  $\ell' \in [L]$  there are  $k'$  edge-dependent message passing layers in  $[\ell']$ , and let  $\ell^*$  be the last such layer. Then  $h_{v_k}^{(\ell^*-1)} = h_{u_k}^{(\ell^*-1)}$  for all  $k > k' - 1$  by the inductive hypothesis. For all  $k > k'$ ,

$$h_{v_k}^{(\ell^*)} = f^{\text{up},(\ell^*)} \left( \min \{ f^{\text{agg},(\ell^*)}(h_w^{(\ell^*-1)} \oplus x_{(w,v_k)}) : w \in \{v_{k-1}, u_{k-1}, v_k, v_{k+1}, u_{k+1}\} \} \oplus h_{v_k}^{(\ell^*-1)} \right)$$

and

$$h_{u_k}^{(\ell^*)} = f^{\text{up},(\ell^*)} \left( \min \{ f^{\text{agg},(\ell^*)}(h_w^{(\ell^*-1)} \oplus x_{(w,u_k)}) : w \in \{v_{k-1}, u_{k-1}, u_k, v_{k+1}, u_{k+1}\} \} \oplus h_{u_k}^{(\ell^*-1)} \right).$$

Recall that  $x_{(w,w)}$  is always zero for any node  $w$ . The minimums in these expressions are then taken over the sets

$$S_v = \{ h_{v_{k-1}}^{(\ell^*-1)} \oplus x_{(v_{k-1}, v_k)}, h_{u_{k-1}}^{(\ell^*-1)} \oplus x_{(u_{k-1}, v_k)}, h_{v_k}^{(\ell^*-1)} \oplus 0, h_{v_{k+1}}^{(\ell^*-1)} \oplus x_{(v_{k+1}, v_k)}, h_{u_{k+1}}^{(\ell^*-1)} \oplus x_{(u_{k+1}, v_k)} \}$$

and

$$S_u = \{ h_{v_{k-1}}^{(\ell^*-1)} \oplus x_{(v_{k-1}, u_k)}, h_{u_{k-1}}^{(\ell^*-1)} \oplus x_{(u_{k-1}, u_k)}, h_{u_k}^{(\ell^*-1)} \oplus 0, h_{v_{k+1}}^{(\ell^*-1)} \oplus x_{(v_{k+1}, u_k)}, h_{u_{k+1}}^{(\ell^*-1)} \oplus x_{(u_{k+1}, u_k)} \}$$

respectively. Note, by the inductive hypothesis,  $h_{v_{k-1}}^{(\ell^*-1)} = h_{u_{k-1}}^{(\ell^*-1)}$ ,  $h_{v_k}^{(\ell^*-1)} = h_{u_k}^{(\ell^*-1)}$ , and  $h_{v_{k+1}}^{(\ell^*-1)} = h_{u_{k+1}}^{(\ell^*-1)}$ , since  $k - 1 > k' - 1$ . Furthermore, for all  $i \in [K]$ , we have  $x_{(v_{i-1}, u_i)} = x_{(u_{i-1}, v_i)} = 1$  and  $x_{(v_{i-1}, v_i)} = x_{(u_{i-1}, u_i)} = 0$ . Thus,

$$S_v = S_u = \{ h_{u_{k-1}}^{(\ell^*-1)} \oplus 0, h_{u_{k-1}}^{(\ell^*-1)} \oplus 1, h_{u_k}^{(\ell^*-1)} \oplus 0, h_{u_{k+1}}^{(\ell^*-1)} \oplus 0, h_{u_{k+1}}^{(\ell^*-1)} \oplus 1 \}.$$

so  $h_{u_k}^{(\ell^*)} = h_{v_k}^{(\ell^*)}$  for  $k > k'$ . Finally, the remaining layers  $\ell \in \{\ell^* + 1, \dots, \ell'\}$  are not edge-dependent message passing, so claim (\*) gives that  $h_{u_k}^{(\ell)} = h_{v_k}^{(\ell)}$  is maintained for  $k > k'$ , completing the induction.

Taking the inductive hypothesis that we just proved with  $\ell' = L$  and  $k' < K$  gives that if there are less than  $K$  edge-dependent message passing layers then  $h_{u_K}^{(L)} = h_{v_K}^{(L)}$ . ■

Note  $x_{v_K}(H_K^{(K)}) \neq x_{u_K}(H_K^{(K)})$ , so any MinAgg GNN with less than  $K$  message passing layers can not achieve perfect accuracy on  $H_K^{(0)}$ . In particular, since a MinAgg GNN with less than  $K$  layers must also have less than  $K$  message passing layers, this lemma demonstrates that  $K$  layers are indeed necessary for a MinAgg GNN to have perfect accuracy on  $K$ -step shortest path problems.

**Corollary 22** *Any MinAgg GNN with less than  $K$  layers has non-zero error on the training pair  $(H_K^{(0)}, H_K^{(K)})$ .*

**Lemma 23** *An  $L$ -layer GNN  $\mathcal{A}_\theta$  has output satisfying*

$$h_{v_0}^{(L)}(P_1^{(0)}(1)) \neq h_{v_1}^{(L)}(P_1^{(0)}(1)) \quad (16)$$

only if for all  $\ell \in [L]$  the function  $f^{\text{up},(\ell)}$  is not constant.

**Proof** Suppose  $f^{\text{up},(\ell)}$  is constant for some  $\ell \in [L]$ . Then the node features after the  $\ell$ th layer are identical:  $h_{v_0}^{(\ell)} = h_{v_1}^{(\ell)}$ . Consider the attributed graph at this layer  $G^{(\ell)} = (V, E, X_v^{(\ell)}, X_e)$  with node features  $x_v = h_v^{(\ell)}$ . This graph has

$$\begin{aligned} \phi : v_0 &\mapsto v_1 \\ v_1 &\mapsto v_0 \end{aligned}$$

as an automorphism. By Fact 17,  $A_\theta(G)$  must also have  $\phi$  as an automorphism. Thus, since isomorphisms respect node features (Def. 16),  $h_{v_0}^{(L)} = h_{\phi(v_0)}^{(L)} = h_{v_1}^{(L)}$ . ■

We are now ready to show that any MinAgg GNN achieving small loss on a training set that includes  $P_1^{(0)}(1)$  and  $H_K^{(0)}$  must have a specific sparsity structure.

**Lemma 24** *For  $K \geq 1$ , let  $\mathcal{G}_{\text{train}}$  be a set containing pairs of training instances  $(G^{(t)}, \Gamma^K(G^{(t)}))$  where  $\mathcal{G}_{\text{train}}$  contains  $M$  total reachable nodes and*

$$\{(H_K^{(0)}, H_K^{(K)}), (P_1^{(0)}(1), P_1^{(K)}(1))\} \subset \mathcal{G}_{\text{train}}.$$

For  $L \geq K > 0$ , consider an  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$  with  $m$ -layer MLPs and parameters  $\theta$ . Given regularization coefficient  $0 < \eta < \frac{1}{2M(mL+mK+K)}$  and error  $0 \leq \epsilon < \eta$ , then the loss

$$\mathcal{L}_{\text{reg}} = \mathcal{L}_{\text{MAE}}(G_{\text{train}}, \mathcal{A}_\theta) + \eta \|\Theta\|_0 \quad (17)$$

has a minimum value of  $\eta(mL + mK + K)$  and if the loss achieved by  $A_\theta$  is within  $\epsilon$  of this minimum then  $A_\theta$  has exactly  $mL + mK + K$  non-zero parameters,  $K$  message passing layers, and for all layers  $f^{\text{up},(\ell)}$  is non-constant. In particular, for all  $j \in [m]$  and  $\ell \in [L]$

$$\|W_j^{\text{up},(\ell)}\|_0 = 1,$$

and  $b_j^{\text{up},(\ell)} = b_j^{\text{agg},(\ell)} = 0$ . For each of the  $K$  message passing layers  $\ell_k \in [L]$ ,

$$\begin{aligned} \|W_1^{\text{agg},(\ell)}\|_0 &= 2 \\ \|W_j^{\text{agg},(\ell)}\|_0 &= 1 \quad \text{for } m \geq j > 1 \end{aligned}$$

where the two non-zero entries in  $W_1^{\text{agg},(\ell)}$  share the same row.

**Proof** First, the BF implementation given by Lemma 12 shows there is a choice of parameters that achieves perfect accuracy with  $mL + mK + K$  non-zero values. This implementation achieves a loss of  $\mathcal{L}_{\text{reg}} = \eta(mL + mK + K)$ . Thus,  $A_\theta$  can not have more than  $mL + mK + K$  non-zero values as this would mean a loss at least  $\eta > \epsilon$  greater than the loss achieved by the BF implementation. We now derive the sparsity structure that any MinAgg GNN achieving a loss no less than  $\eta(mL + mK + K) + \epsilon$  must achieve. We show at the end of the proof that  $\eta(mL + mK + K)$  is indeed the global minimum of the loss.

Note that  $|x_{v_0}(P_1^{(K)}(1)) - x_{v_1}(P_1^{(K)}(1))| = 1$  and  $|x_{v_K}(H_K^{(K)}) - x_{u_K}(H_K^{(K)})| = 1$ . If  $h_{v_0}^{(L)}(P_1^{(0)}(1)) = h_{v_1}^{(L)}(P_1^{(0)}(1))$  then

$$\begin{aligned} \mathcal{L}_{\text{reg}} &\geq \mathcal{L}_{\text{MAE}}(\mathcal{G}_{\text{train}}, A_\theta) \\ &\geq \frac{|x_{v_0}(P_1^{(K)}(1)) - h_{v_0}^{(L)}(P_1^{(0)}(1))| + |x_{v_1}(P_1^{(K)}(1)) - h_{v_1}^{(L)}(P_1^{(0)}(1))|}{M} \\ &\geq \frac{1}{M} \\ &\geq \eta(mL + mK + K) + \epsilon, \end{aligned}$$

where the last inequality follows from  $\epsilon < \frac{1}{2M(mL+mK+K)}$  and  $\eta(mL + mK + K) < \frac{1}{2M}$ . We can then conclude that any MinAgg GNN with loss less than or equal to  $\eta(mL + mK + K) + \epsilon$  must satisfy  $h_{v_k}^{(L)}(H_K^{(0)}) \neq h_{u_k}^{(L)}(H_K^{(0)})$  and  $h_{v_0}^{(L)}(P_1^{(0)}(1)) \neq h_{v_1}^{(L)}(P_1^{(0)}(1))$  so by lemmas 21 and 23, any such MinAgg GNN has at least  $K$  edge-dependent message passing layers and for all  $\ell \in [L]$  the update function  $f^{\text{up},(\ell)}$  is nonconstant.

Next, we analyze the sparsity required to achieve  $K$  edge-dependent message passing layers and nonconstant  $f^{\text{up},(\ell)}$ .

- If  $f^{\text{up},(\ell)}$  is non constant, then for all  $j \in [m]$  and  $\ell \in [L]$

$$\|W_j^{\text{up},(\ell)}\|_0 \geq 1 \tag{18}$$

so there must be at least  $mL$  non-zero entries coming from  $f^{\text{up},(\ell)}$ .

- If  $\ell$  is a edge-dependent message passing layer then  $f^{\text{agg},(\ell)}$  is nonconstant and so for all  $j \in [m]$  and  $\ell \in [L]$

$$\|W_j^{\text{agg},(\ell)}\|_0 \geq 1. \quad (19)$$

Also, since  $f^{\text{agg},(\ell)}$  is edge-dependent message passing it depends on both its node and its edge component so  $W_1^{\text{agg},(\ell)}$  must have two columns that have non-zero entries, giving

$$\|W_1^{\text{agg},(\ell)}\|_0 \geq 2. \quad (20)$$

Thus, the total number of non-zero entries coming from  $f^{\text{agg},(\ell)}$  is at least  $mK + K$  since there are at least  $K$  edge-dependent message passing layers.

Combining the contributions from  $f^{\text{up},(\ell)}$  and  $f^{\text{agg},(\ell)}$  we have that  $A_\theta$  has at least  $mL + mK + K$  non-zero parameters. However, at the start of the proof we showed that  $mL + mK + K$  is also an upper bound on the number of non-zero parameters, so there must be exactly  $mL + mK + K$  non-zero parameters. Additionally, we can see that there can not be more than  $K$  message passing layers as this would require additional non-zero parameters. Furthermore, if only  $mL + mK + K$  parameters are non-zero then inequalities (18) - (20) must be tight and, also,  $b_j^{\text{up},(\ell)} = b_j^{\text{agg},(\ell)} = 0$ . Finally, we remark that the non-zero entries in  $W_1^{\text{agg},(\ell)}$  must share a row since  $W_2^{\text{agg},(\ell)}$  has only one non-zero entry. Thus, if the non-zero entries in  $W_1^{\text{agg},(\ell)}$  do not share a row, either edge dependence or node dependence is lost.

Furthermore, any MinAgg GNN that achieves a loss less than or equal to  $\eta(mL + mK + K) + \epsilon$  must have exactly  $mL + mK + K$  non-zero parameters. This implies that  $\eta(mL + mK + K)$  is the minimum of the loss. ■

We can now limit our analysis to MinAgg GNNs with exactly  $K$  message passing layers. For these MinAgg GNNs let the  $k$ th message passing layer be  $\ell_k \in [L]$  where  $k \in [K]$ .

#### C.4. Bounding GNN expressivity

The following section aims to bound the expressiveness of  $\mathcal{A}_\theta$  under certain conditions that are more challenging to analyze. The benefit is that with sufficiently many training examples we can restrict our analysis to more straightforward cases, as for most training instances, these challenging to analyze conditions do not arise.

**Lemma 25** Consider  $P_{K+1}^{(1)}(x, a_2, \dots, a_{K+1})$  where  $a_2, \dots, a_{K+1}$  are taken to be fixed and view the computed node feature  $h_{v_K}^{(L)} : \mathbb{R} \rightarrow \mathbb{R}$  solely as a function of the first edge weight  $x$ , which is variable. Consider a MinAgg GNN  $\mathcal{A}_\theta$  with exactly  $K$  message passing layers. If for the  $k$ th message passing layer  $\ell_k \in [L]$  there is a region  $D_k \subset \mathbb{R}$  such that the feature  $h_{v_{K+1}}^{(\ell_k)}(x)$  is constant on  $D_k$ , i.e., takes the same value for all  $x \in D_k$ , then  $h_{v_{K+1}}^{(\ell_K)}(x)$  is also constant over  $D_k$ .

**Proof** We proceed by induction, so consider  $k' \in [K]$  with  $k' \geq k$  and assume that  $h_{v_{k'}}^{(\ell_{k'-1})}(x)$  is constant over  $D_k$ . We aim to show  $h_{v_{k'+1}}^{(\ell_{k'})}(x)$  is constant over  $D_k$ . This is true for the base case of

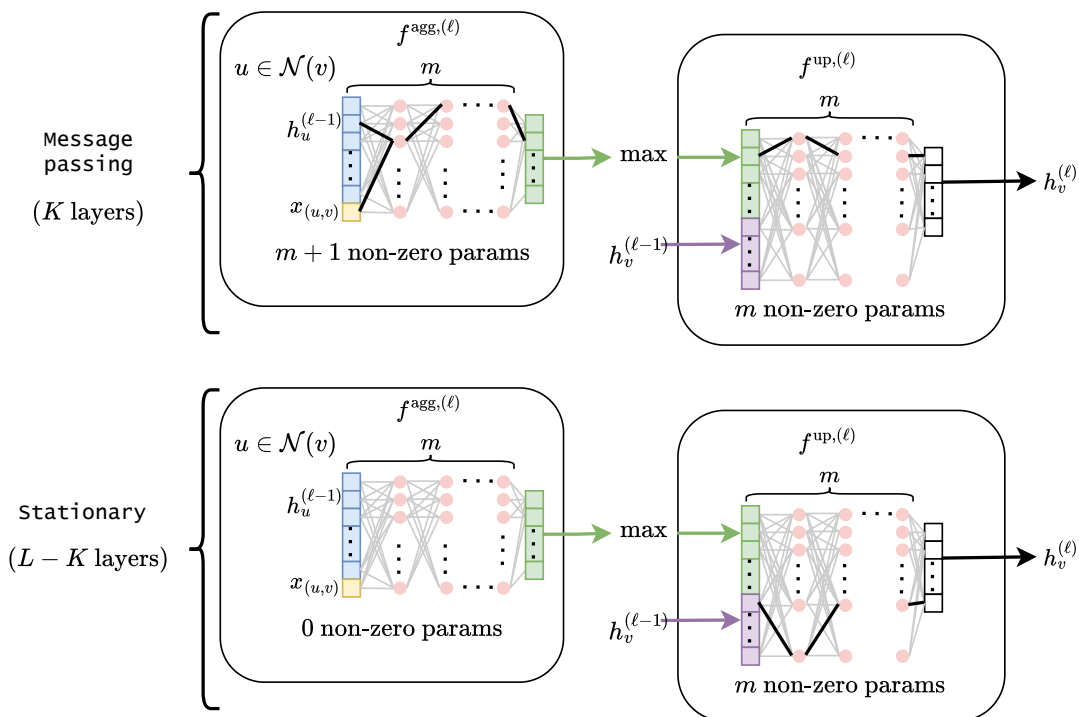


Figure 6: A diagram showing an example of a MinAgg GNN with the sparsity structure given by Lemma 24. Bold black connections in the neural network indicate non zero parameters, while grey lines indicate zero parameters.

$k' = k$  by the assumption in the theorem statement that  $h_{v_{k+1}}^{(\ell_k)}(x)$  is constant on  $D_k$ . Now we prove the general case, so take  $k' > k$ . Since,

$$h_{v_{k'+1}}^{(\ell_{k'})} = f^{\text{up},(\ell_{k'})} \left( \min\{f^{\text{agg},(\ell_{k'})}(h_u^{(\ell_{k'}-1)} \oplus x_{(u,v)}) : u \in \{v_{k'}, v_{k'+1}, v_{k'+2}\}\} \oplus h_{v_{k'+1}}^{(\ell_{k'}-1)} \right), \quad (21)$$

$h_{v_{k'+1}}^{(\ell_{k'})}(x)$  is a function solely of  $h_{v_{k'+1}}^{(\ell_{k'}-1)}(x)$ ,  $h_{v_{k'+2}}^{(\ell_{k'}-1)}(x)$ , and  $h_{v_{k'}}^{(\ell_{k'}-1)}(x)$ . (All edge weights other than that of  $(v_0, v_1)$  are taken to be constant, and  $k' > k \geq 1$  so there is no dependence on  $x_{(v_0, v_1)}$  in the expression.) The feature  $h_{v_{k'}}^{(\ell_{k'}-1)}(x)$  is constant over  $D_k$  by the inductive hypothesis. Furthermore,  $h_{v_{k'}}^{(\ell_{k'}-1)}(x)$  is constant over  $D_k$  by Fact 19 since all layers in  $\{\ell_{k'-1} + 1, \dots, \ell_{k'} - 1\}$  must be stationary. Note that while the stationary layers may change the node feature at  $v_{k'}$  so that  $h_{v_{k'}}^{(\ell_{k'}-1)}(x) \neq h_{v_{k'}}^{(\ell_{k'}-1)}(x)$ , these two features  $h_{v_{k'}}^{(\ell_{k'}-1)}(x)$  and  $h_{v_{k'}}^{(\ell_{k'}-1)}$  are still both constant functions of  $x$  on  $D_k$ . In  $P_{K+1}^{(1)}(x, a_2, \dots, a_{K+1})$  only the first edge weight  $x_{(v_0, v_1)}$  and the second node feature  $x_{v_1}$  depend on  $x$ . Thus, by Claim 20, if  $j \in [K]$  with  $j > k'$  then  $h_{v_{j+1}}^{(\ell_{k'})}(x)$  is constant across all  $x \geq 0$  as  $v_{j+1}$  is more than  $k'$  steps from  $v_1$ . It then follows that the features  $h_{v_{k'+1}}^{(\ell_{k'}-1)}(x)$  and  $h_{v_{k'+2}}^{(\ell_{k'}-1)}(x)$  are both constant. Also,  $h_{v_{k'+1}}^{(\ell_{k'}-1)}(x)$  and  $h_{v_{k'+2}}^{(\ell_{k'}-1)}(x)$  are both constant since the layers in  $\{\ell_{k'-1} + 1, \dots, \ell_{k'} - 1\}$  must be stationary. We can then conclude  $h_{v_{k'+1}}^{(\ell_{k'})}(x)$  is constant on  $D_k$ , since it depends only on variables which are constant on  $D_k$ , completing the inductive argument. Since we have proved  $h_{v_{k'+1}}^{(\ell_{k'})}(x)$  is constant on  $D_k$  for all  $k' \geq k$ , we have that  $h_{v_{K+1}}^{(\ell_K)}(x)$  is constant. Finally, since  $h_{v_{K+1}}^{(m)}(x)$  only depends on  $h_{v_K}^{(\ell_K)}(x)$ , as all layers following  $\ell_k$  are stationary, we have that  $h_{v_K}^{(m)}(x)$  is constant on  $D_k$  ■

**Definition 26** A GNN  $\mathcal{A}_\theta$  has 1-dimensional aggregation if for all message passing steps  $\ell_k \in [L]$

$$\|W_0^{\text{up},(\ell_k)}\|_0 = 1.$$

One dimensional aggregation implies that the output of  $f^{\text{up},(\ell)}$  only depends on one component of the vector produced by min. Furthermore, the value of this one component is determined by which vector in the set

$$\{f^{\text{agg},(\ell_k)}(h_u^{(\ell_k)}, x_{(u,v)}) \mid u \in \mathcal{N}(v)\} \quad (22)$$

is minimal in that component.

**Definition 27** Consider  $P_{K+1}^{(1)}(a_1, a_2, \dots, a_{K+1})$  and a MinAgg GNN  $\mathcal{A}_\theta$  with exactly  $K$  message passing layers and 1-dimensional aggregation. We say that  $\mathcal{A}_\theta$  is path derived on  $P_{K+1}^{(1)}(a_1, a_2, \dots, a_{K+1})$  if for the  $k$ th message passing step  $\ell_k$ ,

$$\begin{aligned} h_{v_{k+1}}^{(\ell_k)} &:= f^{\text{up},(\ell_k)} \left( \min\{f^{\text{agg},(\ell_k)}(h_u^{(\ell_k-1)}, x_{(u,v_{k+1})}) \mid u \in \{v_{k+1}\} \cup N(v_{k+1})\} \oplus h_{v_{k+1}}^{(\ell_k-1)} \right) \\ &= f^{\text{up},(\ell_k)} \left( f^{\text{agg},(\ell_k)}(h_{v_k}^{(\ell_k-1)}, x_{(v_k, v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_k-1)} \right). \end{aligned}$$

**Corollary 28** Consider  $P_{K+1}^{(1)}(x, a_2, \dots, a_{K+1})$  where  $a_2, \dots, a_{K+1}$  are taken to be fixed and view the computed node feature  $h_{v_{K+1}}^{(L)} : \mathbb{R} \rightarrow \mathbb{R}$  solely as a function of the first edge weight  $x$ , which is variable. Consider a MinAgg GNN  $\mathcal{A}_\theta$  with exactly  $K$  message passing layers and 1-dimensional aggregation. If  $D \subseteq \mathbb{R}$  contains all  $x$  such that  $\mathcal{A}_\theta$  is not path derived on  $P_{K+1}^{(1)}(x, a_2, \dots, a_{K+1})$  and

$$Y_{\text{pd}} = \{h_{v_{K+1}}^{(L)}(x) : x \in D\} \quad (23)$$

then  $|Y_{\text{pd}}| \leq K$ .

**Proof**

For  $k \in [K]$ , consider the subset  $D_k \subseteq \mathbb{R}$  such that  $x \in D_k$  if and only if

$$\begin{aligned} & f^{\text{up},(\ell_k)} \left( \min\{f^{\text{agg},(\ell_k)}(h_u^{(\ell_k-1)}, x_{(u,v_{k+1})}) \mid u \in \{v_{k+1}\} \cup N(v_{k+1})\} \oplus h_{v_{k+1}}^{(\ell_k-1)} \right) \\ & \neq f^{\text{up},(\ell_k)} \left( f^{\text{agg},(\ell_k)}(h_{v_{k+1}}^{(\ell_k-1)}, x_{(v_k,v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_k-1)} \right). \end{aligned}$$

Since  $\mathcal{A}_\theta$  has 1-dimensional aggregation, for  $x \in D_k$ ,

$$h_{v_{k+1}}^{(\ell_k)} = f^{\text{up},(\ell_k)} \left( \min\{f^{\text{agg},(\ell_k)}(h_{v_{k+1}}^{(\ell_k-1)}, x_{(v_{k+1},v_{k+1})}), f^{\text{agg},(\ell_k)}(h_{v_{k+2}}^{(\ell_k-1)}, x_{(v_{k+2},v_{k+1})})\} \oplus h_{v_{k+1}}^{(\ell_k-1)} \right)$$

and we have either

$$h_{v_{k+1}}^{(\ell_k)} = f^{\text{up},(\ell_k)} \left( f^{\text{agg},(\ell_k)}(h_{v_{k+1}}^{(\ell_k-1)}, x_{(v_{k+1},v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_k-1)} \right)$$

or,

$$h_{v_{k+1}}^{(\ell_k)} = f^{\text{up},(\ell_k)} \left( f^{\text{agg},(\ell_k)}(h_{v_{k+2}}^{(\ell_k-1)}, x_{(v_{k+2},v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_k-1)} \right).$$

since there is only one component of  $f^{\text{agg},(\ell_k)}(h_{v_{k+1}}^{(\ell_k-1)}, x_{(v_{k+1},v_{k+1})})$  and  $f^{\text{agg},(\ell_k)}(h_{v_{k+2}}^{(\ell_k-1)}, x_{(v_{k+2},v_{k+1})})$  that is not identically zero.

By Claim 20,  $h_{v_{k+1}}^{(\ell_k-1)}(x)$  and  $h_{v_{k+2}}^{(\ell_k-1)}(x)$  are constant functions of  $x$  ( $v_{k+1}$  and  $v_{k+2}$  are both more than  $k-1$  steps away from  $v_1$  and  $(v_0, v_1)$ , and at step  $\ell_k-1$  only  $k-1$  message passing steps have occurred). Then,  $h_{v_{k+1}}^{(\ell_k)}(x)$  must be constant over  $x \in D_k$  since it only depends on features that are constant. By Lemma 25,  $h_{v_{K+1}}^{(L)}(x)$  is constant on  $D_k$  and must take some value  $q_k$ . Let  $Q = \{q_k \mid k \in [K]\}$  and note  $|Q| \leq K$ .

Suppose  $h_{v_{K+1}}^{(L)}(x)$  is not path derived. Then there exists some  $k \in [K]$  such that

$$\begin{aligned} & f^{\text{up},(\ell_k)} \left( \min\{f^{\text{agg},(\ell_k)}(h_u^{(\ell_k-1)}, x_{(u,v_{k+1})}) \mid u \in \{v_{k+1}\} \cup N(v_{k+1})\} \oplus h_{v_{k+1}}^{(\ell_k-1)} \right) \\ & \neq f^{\text{up},(\ell_k)} \left( f^{\text{agg},(\ell_k)}(h_{v_{k+1}}^{(\ell_k-1)}, x_{(v_k,v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_k-1)} \right). \end{aligned}$$

which implies  $h_{v_{K+1}}^{(L)}(x) \in Q$  so  $Y_{\text{pd}} \subset Q$  and  $|Y_{\text{pd}}| \leq K$ . ■

### C.5. Global minimum is BF

We are now ready to prove Theorem 3, which we restate here with additional details.

**Theorem 29** *Let  $\mathcal{G}_{\text{train}}$  be a set containing pairs of training instances  $(G^{(t)}, \Gamma^K(G^{(t)}))$  where  $\mathcal{G}_{\text{train}}$  contains  $M$  total reachable nodes and  $\mathcal{G}_K \subset \mathcal{G}_{\text{train}}$ . For  $L \geq K > 0$ , consider an  $L$ -layer MinAgg GNN  $\mathcal{A}_\theta$  with  $m$ -layer MLPs and parameters  $\theta$ . Given regularization coefficient  $0 < \eta < \frac{1}{2M(mL+mK+K)}$  and error  $0 \leq \epsilon < \eta$ , then the loss*

$$\mathcal{L}_{\text{reg}} = \mathcal{L}_{\text{MAE}}(\mathcal{G}_{\text{train}}, \mathcal{A}_\theta) + \eta \|\Theta\|_0 \quad (24)$$

*has a minimum value of  $\eta(mL + mK + K)$  and if the loss achieved by  $\mathcal{A}_\theta$  is within  $\epsilon$  of this minimum then on any  $G \in \mathcal{G}$  the features computed by the MinAgg GNN satisfy*

$$(1 - M\epsilon)x_v(\Gamma^K(G)) \leq h_v^{(L)}(G) \leq (1 + M\epsilon)x_v(\Gamma^K(G))$$

*for all  $v \in V(G)$ .*

**Proof** Our proof proceeds by first simplifying and reparametrizing the MinAgg GNN update using the sparsity structure derived Section C.3. Next, we prove approximations to these new parameters, where we utilize the results of Section C.4 to restrict analysis to cases where computation follows a simple structure (path-derived cases). We conclude by using these approximations to bound the error the MinAgg GNN achieves on an arbitrary graph. Key to this final argument is showing that the features of the MinAgg GNN approximate the node values in the BF algorithm, up to some scaling factor.

**Simplifying the update.** We next show that the MinAgg GNN under the sparsity constraints of Lemma 24 can be reformulated into a much simpler framework, where all the information exchanged between message-passing layers is consolidated into a much smaller set of parameters and where node features are always one dimensional. As  $\{(H_K^{(0)}, H_K^{(K)}), (P_1^{(0)}(1), P_1^{(K)}(1))\} \subset \mathcal{G}_K \subset \mathcal{G}_{\text{train}}$ , Lemma 24 gives that  $\mathcal{A}_\theta$  has exactly  $mL + mK + K$  non-zero parameters,  $K$  edge-dependent message passing layers, and a specific sparsity structure: for all layers  $\ell \in [L]$ , and  $j \in [m]$

$$\|W_j^{\text{up},(\ell)}\|_0 = 1$$

and for each of the  $K$  message passing layers  $\ell_k \in [L]$

$$\begin{aligned} \|W_1^{\text{agg},(\ell_k)}\|_0 &= 2 \\ \|W_j^{\text{agg},(\ell_k)}\|_0 &= 1 \quad \text{for } j > 1 \end{aligned}$$

where the two non-zero entries in  $W_1^{\text{agg},(\ell_k)}$  must share a row. Furthermore, all bias terms are zero.

We now argue that each of these non-zero parameters must also be non-negative. For  $\ell \in [L]$ , suppose that there is some  $W_j^{\text{up},(\ell)}$  with a negative element. Then  $\sigma(W_j^{\text{up},(\ell)}x) = 0$ , where  $x$  is a vector with non-negative entries, is a constant function of  $x$ . This implies  $f^{\text{up},(\ell)}$  is constant (it is the zero function), which contradicts Lemma 24. Similarly, for any message passing layer  $\ell_k \in [L]$  if there is some  $W_j^{\text{agg},(\ell_k)}$  with a negative element for  $j > 1$  then  $\sigma(W_j^{\text{agg},(\ell_k)}x) = 0$  is a constant

function of  $x$  which gives that  $f^{\text{agg},(\ell_k)}$  is constant. Again, we have a contradiction with Lemma 24, which says that  $f^{\text{agg},(\ell_k)}$  must be edge and node dependent.

It remains to show that for messages passing layers  $\ell_k \in [L]$  that the non-zero elements in  $W_1^{\text{agg},(\ell_k)}$  are non-negative. However, we first simplify the update function given the restrictions we have already derived. By the sparsity of the MinAgg GNN ( $\|W_m^{\text{up},(\ell)}\|_0 = 1$ ), at any step  $\ell$ , all node features  $h_v^{(\ell)}$  have at most 1 non-zero component. We can then simplify the update by reducing the number of dimensions. In particular we get

$$\tilde{h}_v^{(\ell)} = \gamma^{(\ell)} \tilde{h}_v^{(\ell-1)}$$

for stationary layers and

$$\tilde{h}_v^{(\ell)} = \gamma^{(\ell)} \min\{\sigma(\rho^{(\ell)} \tilde{h}_v^{(\ell-1)} + \tau^{(\ell)} x_{(u,v)}) : u \in \mathcal{N}(v)\}$$

for message passing layer where  $\gamma^{(\ell)}, \rho^{(\ell)}, \tau^{(\ell)} \in \mathbb{R}$  are new parameters that are functions the initial parameters  $\theta$ , and  $\tilde{h}_v^{(\ell)} \in \mathbb{R}$  is a single-dimensional node feature that is equal to the non-zero value in  $h_v^{(\ell)}$  (or zero if there is no such value). We have that  $\tilde{h}_v^{(0)} = h_v^{(0)} = x_v^{(0)}$  and  $\tilde{h}_v^{(L)} = h_v^{(L)}$  since  $h_v^{(0)}, h_v^{(L)} \in \mathbb{R}$ . Note that the new parameters do not depend on the input features since they are only dependent on  $\theta$ .

We now further simplify by combining the stationary layers with their succeeding message passing layer. Let  $\ell_0 = 0$ . For the  $k$ th message passing layer  $\ell_k \in [L]$

$$\begin{aligned} \tilde{h}_v^{(\ell_k)} &= \gamma^{(\ell_k)} \min \left\{ \sigma \left( \rho^{(\ell_k)} \left( \prod_{i \in \{\ell_{k-1}+1, \dots, \ell_k-1\}} \gamma^{(i)} \right) \tilde{h}_v^{(\ell_{k-1})} + \tau^{(\ell_k)} x_{(u,v)} \right) : u \in \mathcal{N}(v) \right\} \\ &= \gamma^{(\ell_k)} \min \left\{ \sigma \left( \alpha^{(\ell_k)} \tilde{h}_v^{(\ell-1)} + \tau^{(\ell_k)} x_{(u,v)} \right) : u \in \mathcal{N}(v) \right\} \end{aligned}$$

where  $\alpha^{(\ell_k)} = \rho^{(\ell_k)} \left( \prod_{i \in \{\ell_{k-1}+1, \dots, \ell_k-1\}} \gamma^{(i)} \right)$ . Note that if for any  $\ell_k \in [L]$  we have  $h_v^{(\ell_k)} = 0$ , then or all succeeding layers  $\ell > \ell_k$ , the node feature  $h_v^{(\ell)}$  is also zero. This is because the min aggregation at  $h_v^{(\ell_k)}$  always includes  $h_v^{(\ell_{k-1})}$ , so if  $h_v^{(\ell_{k-1})} = 0$  then  $h_v^{(\ell_k)} = 0$ . We are now ready to show  $\alpha^{(\ell_k)} > 0$  and  $\tau^{(\ell_k)} > 0$ . Suppose  $\alpha^{(\ell_k)} \leq 0$  for some  $\ell_k \in [L]$  and consider the training instance  $(P_2^{(1)}(1, 0), P_2^{(K+1)}(1, 0))$ . The update at  $v_2$  is

$$\begin{aligned} \tilde{h}_{v_2}^{(\ell_k)} &= \gamma^{(\ell_k)} \min \left\{ \sigma \left( \alpha^{(\ell)} \tilde{h}_u^{(\ell_{k-1})} \right) : u \in \{v_1, v_2\} \right\} \\ &= 0 \end{aligned}$$

since  $h_v^{(\ell_{k-1})} \geq 0$  for all  $v \in V$ . We then also have  $\tilde{h}_{v_2}^{(L)} = 0$  and since  $x_{v_2}(P_2^{(K+1)}(1, 0)) = 1$  the loss is

$$\begin{aligned} \mathcal{L}_{\text{reg}} &\geq \mathcal{L}_{\text{MAE}}(\mathcal{G}_{\text{train}}, \mathcal{A}_\theta) \\ &> 1/M \\ &> \eta(mL + mL + K) + \epsilon \end{aligned}$$

which is a contradiction.

Now instead suppose  $\tau^{(\ell_k)} < 0$  for some  $\ell_k \in [L]$  and consider the training instance  $(P_1^{(0)}(1), P_1^{(K)}(1))$ . Since  $h_{v_0}^{(\ell)} = 0$  for all  $\ell \in [L]$ , the update at  $v_1$  is

$$\begin{aligned}\tilde{h}_{v_1}^{(\ell_k)} &= \gamma^{(\ell_k)} \min \left\{ \sigma \left( \tau^{(\ell_k)} x_{(v_0, v_1)} \right), \sigma \left( \alpha^{(\ell)} \tilde{h}_{v_1}^{(\ell_{k-1})} + \tau^{(\ell_k)} x_{(v_1, v_1)} \right) \right\} \\ &= \gamma^{(\ell_k)} \min \left\{ 0, \sigma \left( \alpha^{(\ell)} \tilde{h}_{v_1}^{(\ell_{k-1})} + \tau^{(\ell_k)} x_{(v_1, v_1)} \right) \right\} \\ &= 0.\end{aligned}$$

We then also have  $\tilde{h}_{v_1}^{(L)} = 0$  and since  $x_{v_1}(P_1^{(K)}(1)) = 1$  the loss is again greater than  $\eta(mL + mK + K) + \epsilon$  which is a contradiction.

We can now remove the final ReLU, since its argument is always non-negative:

$$\tilde{h}_v^{(\ell_k)} = \gamma^{(\ell_k)} \min \left\{ \alpha^{(\ell_k)} \tilde{h}_v^{(\ell_{k-1})} + \tau^{(\ell_k)} x_{(u, v)} : u \in \mathcal{N}(v) \right\}.$$

As another simplification, we re-index to  $k \in [K]$  as follows. Let  $\bar{h}_v^{(k)} = \tilde{h}_v^{(\ell_k)}$  and  $\bar{\gamma}^{(k)} = \gamma^{(\ell_k)}$  for  $k \in \{0, \dots, K-1\}$ . However, to account for the effect of the layers succeeding  $\ell_k$  we take  $\bar{h}_v^{(K)} = \tilde{h}_v^{(L)}$  and  $\bar{\gamma}^{(K)} = \prod_{i \in \{\ell_K, \ell_{K+1}, \dots, L\}} \gamma^{(i)}$ . Furthermore, for all  $k \in [K]$ , let  $\bar{\alpha}^{(k)} = \alpha^{(\ell_k)}$  and  $\bar{\tau}^{(k)} = \tau^{(\ell_k)}$ . Then

$$\bar{h}_v^{(k)} = \bar{\gamma}^{(k)} \min \left\{ \bar{\alpha}^{(k)} \bar{h}_v^{(k-1)} + \bar{\tau}^{(k)} x_{(u, v)} : u \in \mathcal{N}(v) \right\}.$$

Finally, by letting  $\mu^{(k)} = \bar{\gamma}^{(k)} / \bar{\alpha}^{(k)}$  and  $\nu^{(k)} = \bar{\tau}^{(k)} / \bar{\alpha}^{(k)}$  we get

$$\bar{h}_v^{(k)} = \mu^{(k)} \min \left\{ \bar{h}_v^{(k-1)} + \nu^{(k)} x_{(u, v)} : u \in \mathcal{N}(v) \right\}. \quad (25)$$

Note that we can factor through the min here because  $\alpha^{(k)} > 0$ . For the rest of the proof, we focus on this simplified update, since it has the same output as the MinAgg GNN, i.e.,  $\bar{h}_v^{(K)} = h_v^{(L)}$ .

**Approximating parameters.** We proceed by analyzing the output of the MinAgg GNN on an inputs for which it is path derived. To this end, suppose that  $\mathcal{A}_\theta$  is path derived on some input graph  $P_{K+1}^{(1)}(a_1, \dots, a_{K+1})$ . Then, by definition of path derived,

$$h_{v_{k+1}}^{(\ell_k)} = f^{\text{up}, (\ell_k)} \left( f^{\text{agg}, (\ell_k)}(h_{v_k}^{(\ell_{k-1})}, x_{(v_k, v_{k+1})}) \oplus h_{v_{k+1}}^{(\ell_{k-1})} \right)$$

which implies

$$\bar{h}_{v_{k+1}}^{(k)} = \mu^{(k)} (\bar{h}_{v_k}^{(k-1)} + \nu^{(k)} x_{(v_k, v_{k+1})}).$$

Next, we prove, for  $k \in \{0, \dots, K\}$ ,

$$\bar{h}_{v_{k+1}}^{(k)} = \left( \prod_{k \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^k \nu^{(s)} \left( \prod_{k \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})}$$

by induction, where the products evaluate to 1 if they are indexed over an empty set. For the base case we have

$$\bar{h}_{v_1}^{(0)} = x_{(v_0, v_1)}.$$

Now for  $k \in [K]$ , suppose

$$\bar{h}_{v_k}^{(k-1)} = \left( \prod_{k-1 \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^{k-1} \nu^{(s)} \left( \prod_{k-1 \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})}.$$

Then we have

$$\begin{aligned} \bar{h}_{v_{k+1}}^{(k)} &= \mu^{(k)} (\bar{h}_{v_k}^{(k-1)} + \nu^{(k)} x_{(v_k, v_{k+1})}) \\ &= \mu^{(k)} \left( \left( \prod_{k-1 \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^{k-1} \nu^{(s)} \left( \prod_{k-1 \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})} + \nu^{(k)} x_{(v_k, v_{k+1})} \right) \\ &= \left( \prod_{k \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^{k-1} \nu^{(s)} \left( \prod_{k \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})} + \nu^{(k)} \mu^{(k)} x_{(v_k, v_{k+1})} \\ &= \left( \prod_{k \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^k \nu^{(s)} \left( \prod_{k \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})}. \end{aligned}$$

This completes the inductions and yields

$$\bar{h}_{v_{K+1}}^{(K)} = \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) x_{(v_0, v_1)} + \sum_{s=1}^K \nu^{(s)} \left( \prod_{K \geq i \geq s} \mu^{(i)} \right) x_{(v_s, v_{s+1})}$$

Given this expression for the output at  $v_{K+1}$  we can now derive bounds on the values of these parameters. However, we must restrict our focus to instances of the training set for which  $A_\theta$  is path derived.

For  $k \in [K]$ , let  $\mathcal{H}_{k,K}^0$  contain graphs in  $\mathcal{H}_{k,K}$  where the  $k+1$  edge has weight 0 and let  $\mathcal{H}_{k,K}^1$  contain graphs in  $\mathcal{H}_{k,K}$  where the  $k+1$  edge has weight  $2K+1$ . It can be checked that for any  $G^{(1)}, G^{(1')} \in \mathcal{H}_{k,K}$ ,

$$|x_{v_{K+1}}(G^{(K+1)}) - x_{v_{K+1}}(G^{(K+1)})| \geq 1.$$

Then, it must be that  $\bar{h}_{v_{K+1}}^{(K)}(G^{(1)}) \neq \bar{h}_{v_{K+1}}^{(K)}(G^{(1)})$  as if these output features are equal

$$\begin{aligned} \mathcal{L}_{\text{reg}} &\geq \mathcal{L}_{\text{MAE}} \\ &\geq \frac{|\bar{h}_{v_{K+1}}^{(K)}(G^{(1)}) - x_{v_{K+1}}(G^{(K+1)})| + |\bar{h}_{v_{K+1}}^{(K)}(G^{(1)}) - x_{v_{K+1}}(G^{(K+1)})|}{M} \\ &\geq \frac{1}{M} \\ &\geq \eta(mL + mK + K) + \epsilon \end{aligned}$$

which violates that  $|\mathcal{L}_{\text{reg}} - \eta(mL + mK + K)| < \epsilon$ .

Consider the subset  $\mathcal{I}_{k,K}^0 \subset \mathcal{H}_{k,K}^0$  containing all graphs  $G \in \mathcal{H}_{k,K}^0$  for which  $A_\theta$  is not path derived on  $G$ . By Lemma 28, if

$$Y_{pd}^0 = \{\bar{h}_{v_{k+1}}^{(k)}(G) : G \in \mathcal{I}_{k,K}^0\}$$

then  $|Y_{pd}^0| \leq K$  and  $|\mathcal{I}_{k,K}^0| \leq K$ . Similarly, if  $\mathcal{I}_{k,K}^1 \subset \mathcal{H}_{k,K}^1$  contains all graphs  $G \in \mathcal{H}_{k,K}^1$  for which  $A_\theta$  is not path derived on  $G$ , and

$$Y_{pd}^1 = \{\bar{h}_{v_{k+1}}^{(k)}(G) : G \in \mathcal{I}_{k,K}^1\}$$

then  $|Y_{pd}^1| \leq K$  and  $|\mathcal{I}_{k,K}^1| \leq K$ .

Using these facts, the pigeon hole principle gives that there must be two graphs  $G_0^{(1)} \in \mathcal{H}_{k,K}^0 \setminus \mathcal{I}_{k,K}^0$  and  $G_1^{(1)} \in \mathcal{H}_{k,K}^1 \setminus \mathcal{I}_{k,K}^1$  such that

$$x_{(v_0,v_1)}(G_0^{(1)}) = x_{(v_0,v_1)}(G_1^{(1)}),$$

i.e.,  $G_0^{(1)}$  and  $G_1^{(1)}$  only differ in the weight of their  $(k+1)$ th edge.

Since error less than  $\epsilon$  is achieved on  $\mathcal{G}_{\text{train}}$ ,

$$|\bar{h}_{v_{K+1}}^{(K)}(G_0^{(1)}) - x_{v_{K+1}}(G_0^{(K+1)})| \leq M\epsilon$$

and

$$|\bar{h}_{v_{K+1}}^{(K)}(G_1^{(1)}) - x_{v_{K+1}}(G_1^{(K+1)})| \leq M\epsilon.$$

The triangle inequity then gives

$$|\bar{h}_{v_{K+1}}^{(K)}(G_0^{(1)}) - x_{v_{K+1}}(G_0^{(K+1)}) - \bar{h}_{v_{K+1}}^{(K)}(G_1^{(1)}) + x_{v_{K+1}}(G_1^{(K+1)})| \leq 2M\epsilon.$$

Making the substitutions

$$\bar{h}_{v_{K+1}}^{(K)}(G_0^{(1)}) = \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) x_{(v_0,v_1)}(G_0^{(1)}) + \sum_{s=1}^K \nu^{(s)} \left( \prod_{K \geq i \geq s} \mu^{(i)} \right) x_{(v_s,v_{s+1})}(G_0^{(1)})$$

$$\bar{h}_{v_{K+1}}^{(K)}(G_1^{(1)}) = \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) x_{(v_0,v_1)}(G_1^{(1)}) + \sum_{s=1}^K \nu^{(s)} \left( \prod_{K \geq i \geq s} \mu^{(i)} \right) x_{(v_s,v_{s+1})}(G_1^{(1)})$$

$$x_{v_{K+1}}(G_0^{(K+1)}) = x_{(v_0,v_1)}(G_0^{(1)}) + x_{(v_k,v_{k+1})}(G_0^{(1)})$$

$$x_{v_{K+1}}(G_1^{(K+1)}) = x_{(v_0,v_1)}(G_1^{(1)}) + x_{(v_k,v_{k+1})}(G_1^{(1)})$$

and canceling like terms yields

$$|\nu^{(k)} \left( \prod_{K \geq i \geq k} \mu^{(i)} \right) (x_{(v_k,v_{k+1})}(G_0^{(1)}) - x_{(v_k,v_{k+1})}(G_1^{(1)}) - (x_{(v_k,v_{k+1})}(G_0^{(1)}) - x_{(v_k,v_{k+1})}(G_1^{(1)}))| \leq 2M\epsilon.$$

Since  $|x_{(v_k, v_{k+1})}(G_0^{(1)}) - x_{(v_k, v_{k+1})}(G_1^{(1)})| = 2K + 1 \geq 1$ , dividing by this factor gives

$$|\nu^{(k)} \left( \prod_{K \geq i \geq k} \mu^{(i)} \right) - 1| \leq 2M\epsilon.$$

We can rewrite the inequality as

$$\frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} (1 - 2M\epsilon) \leq \nu^{(k)} \mu^{(k)} \leq \frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} (1 + 2M\epsilon). \quad (26)$$

Next we bound the product  $\prod_{K \geq i \geq 1} \mu^{(i)}$  by considering how the MinAgg GNN scales the edge weight  $x_{(v_0, v_1)}$ . Consider two instances  $J^{(1)}, J'^{(1)} \in \mathcal{H}_{k, K}^0 \setminus \mathcal{J}_{k, K}^0$ , that is, instances for which  $A_\theta$  is path derived (which must exist since  $|\mathcal{J}_{k, K}^0| \leq K$  and  $|\mathcal{H}_{k, K}^0| = 2K + 1$ ). These instances only differ in the weight of their first edge. As before,

$$|\bar{h}_{v_{K+1}}^{(K)}(J^{(1)}) - x_{v_{K+1}}(J^{(K+1)})| \leq M\epsilon$$

and

$$|\bar{h}_{v_{K+1}}^{(K)}(J'^{(1)}) - x_{v_{K+1}}(J'^{(K+1)})| \leq M\epsilon.$$

which combine to give

$$|\bar{h}_{v_{K+1}}^{(K)}(J^{(1)}) - x_{v_{K+1}}(J^{(K+1)}) - \bar{h}_{v_{K+1}}^{(K)}(J'^{(1)}) + x_{v_{K+1}}(J'^{(K+1)})| \leq 2M\epsilon.$$

Substituting in for these four terms and canceling yields

$$\left| \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) (x_{(v_0, v_1)}(J^{(1)}) - x_{(v_0, v_1)}(J'^{(1)})) - (x_{(v_0, v_1)}(J^{(1)}) - x_{(v_0, v_1)}(J'^{(1)})) \right| \leq 2M\epsilon.$$

Since  $|x_{(v_0, v_1)}(J^{(1)}) - x_{(v_0, v_1)}(J'^{(1)})| \geq 1$  we can rearrange this inequality as

$$\left| \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) - 1 \right| \leq 2M\epsilon.$$

and

$$(1 - 2M\epsilon) \leq \left( \prod_{K \geq i \geq 1} \mu^{(i)} \right) \leq (1 + 2M\epsilon). \quad (27)$$

**Bounding error on arbitrary graphs.** Consider an arbitrary graph  $G \in \mathcal{G}$ . We aim to show that the MinAgg GNN outputs on this graph  $h_v^{(L)} = \bar{h}_v^{(K)}$  approximate  $x_v(\Gamma^K(G))$ . Let  $r_v^{(0)} = x_v(G)$  and define

$$r_v^{(k)} = \min\{r_u^{(k-1)} + x_{(u, v)} \mid u \in \mathcal{N}(v)\}.$$

meaning  $r_v^{(k)}$  are the results of applying  $k$  steps of the BF algorithm. In particular,  $r_v^{(K)}(G) = x_v(\Gamma^K(G))$

We now show that

$$\bar{h}_v^{(k)} \leq \frac{1 + 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} r_v^{(k)}$$

by induction. The base case

$$\bar{h}_v^{(0)} = r_v^{(0)} \leq \frac{1 + 2M\epsilon}{\prod_{K \geq i \geq 1} \mu^{(i)}} r_v^{(0)}$$

follows from Eq. 27. Now suppose

$$\bar{h}_v^{(k-1)} \leq \frac{1 + 2M\epsilon}{\prod_{K \geq i \geq k} \mu^{(i)}} r_v^{(k-1)}.$$

Then, using Eq. 26,

$$\begin{aligned} \bar{h}_v^{(k)} &= \mu^{(k)} \min \left\{ \bar{h}_v^{(k-1)} + \nu^{(k)} x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\ &\leq \mu^{(k)} \min \left\{ (1 + 2M\epsilon) r_v^{(k-1)} \frac{1}{\prod_{K \geq i \geq k} \mu^{(i)}} + \nu^{(k)} x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\ &\leq \min \left\{ (1 + 2M\epsilon) r_v^{(k-1)} \frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} + \frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} (1 + 2M\epsilon) x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\ &= \frac{1 + 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} \min \left\{ r_v^{(k-1)} + x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\ &= \frac{1 + 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} r_v^{(k)}. \end{aligned}$$

On the other hand, we next show

$$\bar{h}_v^{(k)} \geq \frac{1 - 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} r_v^{(k)}$$

by induction. The base case

$$\bar{h}_v^{(0)} = r_v^{(0)} \geq \frac{1 - 2M\epsilon}{\prod_{K \geq i \geq 1} \mu^{(i)}} r_v^{(0)}$$

again follows from Eq. 27. Now suppose

$$\bar{h}_v^{(k-1)} \geq \frac{1 - 2M\epsilon}{\prod_{K \geq i \geq k} \mu^{(i)}} r_v^{(k-1)}.$$

Then, using Eq. 26,

$$\begin{aligned}
 \bar{h}_v^{(k)} &= \mu^{(k)} \min \left\{ \bar{h}_v^{(k-1)} + \nu^{(k)} x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\
 &\geq \mu^{(k)} \min \left\{ (1 - 2M\epsilon) r_v^{(k-1)} \frac{1}{\prod_{K \geq i \geq k} \mu^{(i)}} + \nu^{(k)} x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\
 &\geq \min \left\{ (1 - 2M\epsilon) r_v^{(k-1)} \frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} + \frac{1}{\prod_{K \geq i \geq k+1} \mu^{(i)}} (1 - 2M\epsilon) x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\
 &= \frac{1 - 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} \min \left\{ r_v^{(k-1)} + x_{(u,v)} : u \in \mathcal{N}(v) \right\} \\
 &= \frac{1 - 2M\epsilon}{\prod_{K \geq i \geq k+1} \mu^{(i)}} r_v^{(k)}.
 \end{aligned}$$

Putting the upper bound and lower bound together, taking  $k = K$ , and using  $x_v(\Gamma^K(G)) = r_v^{(K)}$  yields

$$(1 - 2M\epsilon)x_v(\Gamma^K(G)) \leq \bar{h}_v^{(K)} \leq (1 + 2M\epsilon)x_v(\Gamma^K(G))$$

■

## Appendix D. Additional Experiments

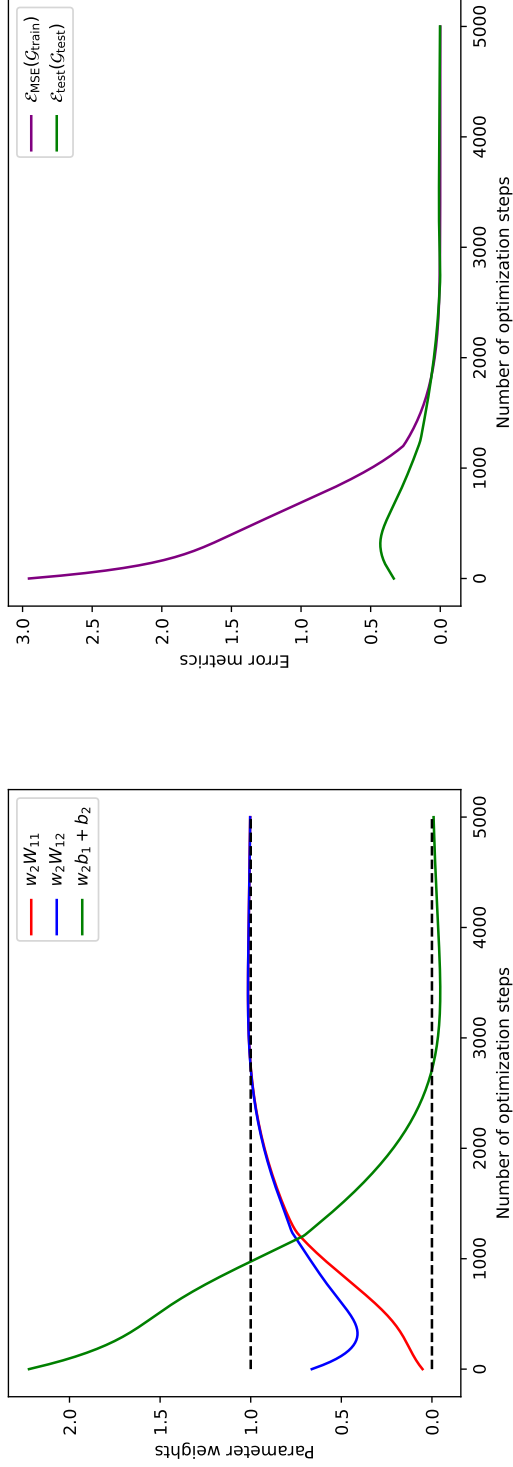
We include additional experiments verifying our theoretical claims with different configurations of the MinAgg GNN including the simple BF-GNN described in Theorem 2 as well as several configurations of the complex BF-GNN utilized in Theorem 3. All models are trained on 8 NVidia A100 GPUs using the AdamW optimizer with a learning rate of 0.001. We evaluate each model using the metrics described in Section 3.

### D.1. Simple MinAgg GNN

We show in Fig. 7 that, with gradient descent, a simple MinAgg GNN will converge to the parameter configuration described in Theorem 2. The simple MinAgg GNN is trained with the specified train set in Theorem 2: four single-edge graphs initialized at step 0 and two double-edge graphs initialized at step 1. Recall that an update for the simple MinAgg GNN is defined as

$$h_v^{(1)} = \sigma(w_2 \min_{u \in \mathcal{N}(v)} \{\sigma(W_1(x_u \oplus x_{(v,u)} + b_1))\} + b_2)$$

where  $w_2, b_1, b_2 \in \mathbb{R}$  and  $W_1 \in \mathbb{R}^{1 \times 2}$ . From Theorem 2, we know that the simple BF-GNN will implement a single step of Bellman-Ford if  $w_2 W_{11} = w_2 w_{12} = 1$  and  $w_2 b_1 + b_2 = 0$ . Therefore, in Fig. 7 (a), we see that, via gradient descent, the parameter configurations converge to the expected values (indicated by the black dotted lines). We further empirically verify the results in Theorem 2 in Fig. 7(b) by showing that as the parameter configurations converge to the expected values, the test error also converges to zero.



(a) Model weights for small BF-GNN.

(b) Error metrics for small BF-GNN.

Figure 7: Performance and parameter weights for the small BF-GNN instance (trained with MSE). Recall from to Theorem 2, we expect that  $w_2 b_1 + b_2 = 0$  and  $w_2 W_{11} = w_2 W_{12} = 1.0$ . We indicate these values by the black dotted lines in (a) and show that the parameters of our small BF-GNN converge to the expected parameter values from Theorem 2. Additionally, we verify that convergence to this parameter configuration corresponds to low test error in (b) as we have that  $\mathcal{E}_{\text{approx}}$  converges to 0.0018.

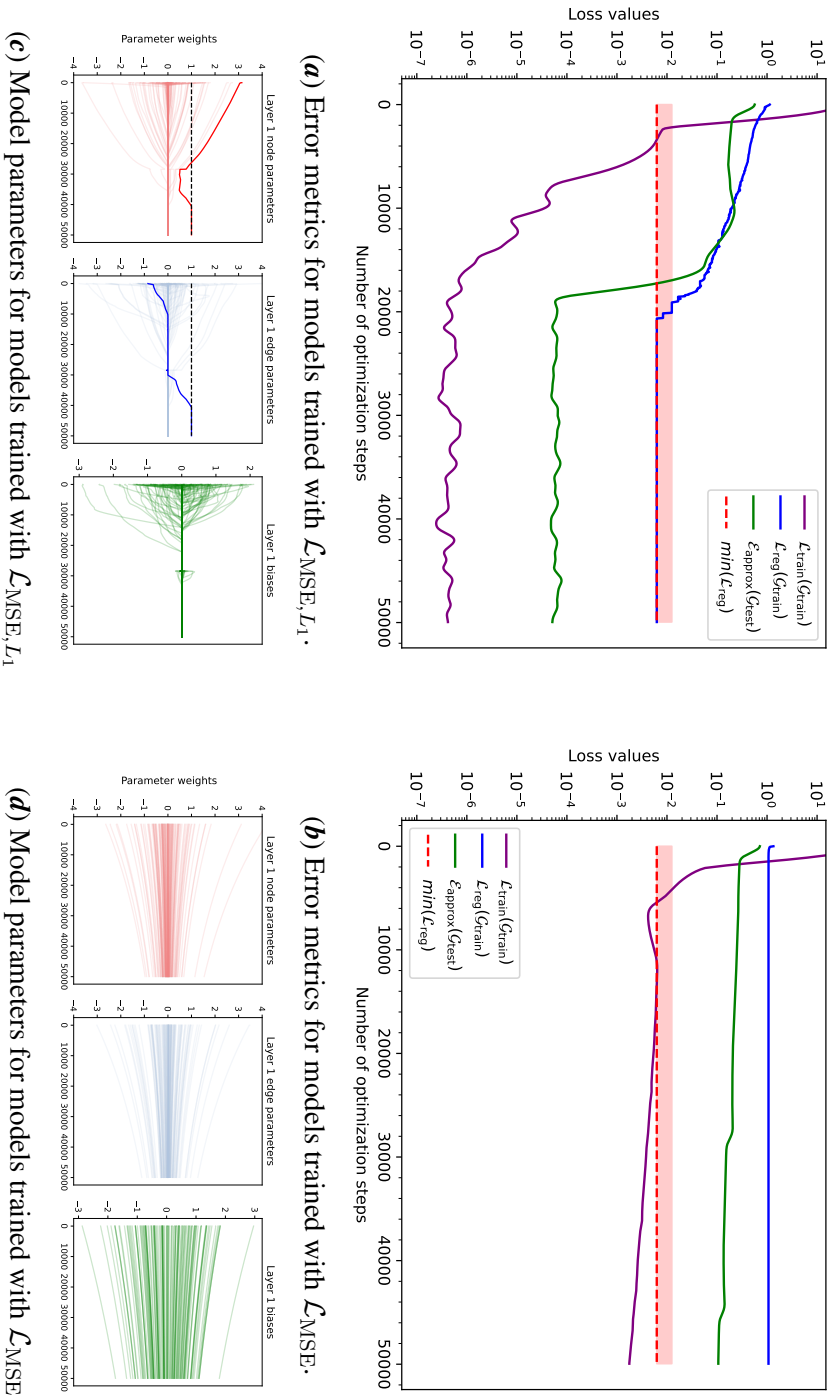
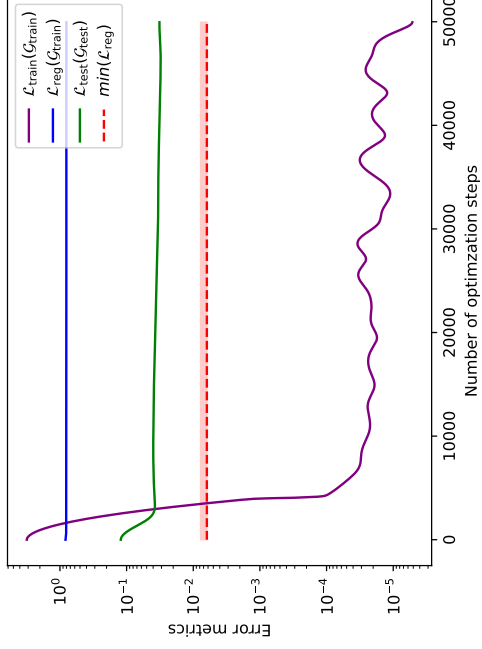
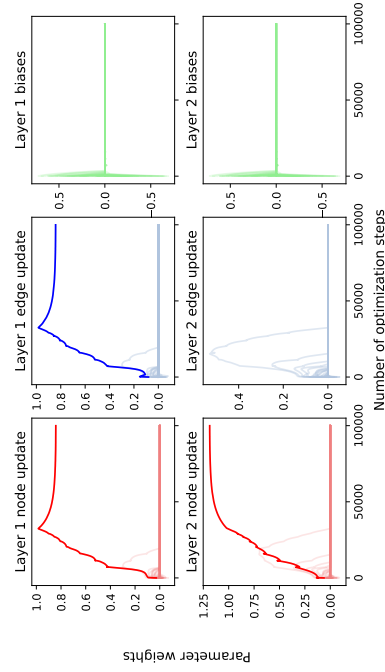


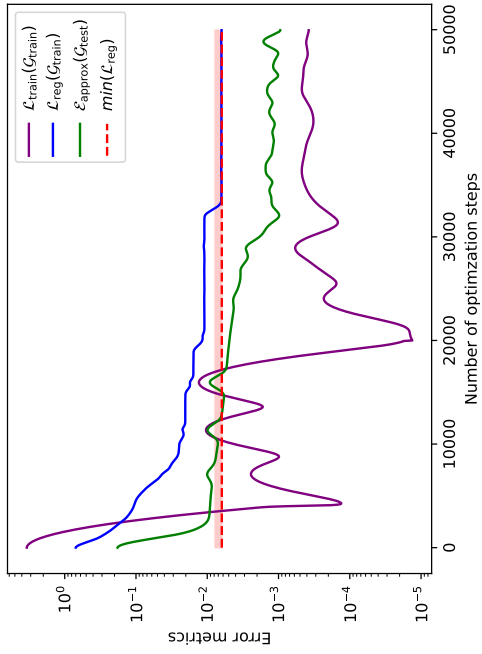
Figure 8: Error metrics and parameter updates for a one-layer MiniAgg GNN trained on a single step of the Bellman-Ford algorithm. Note that the dotted line in (a) and (b) is the global minimum of  $\mathcal{L}_{\text{reg}}$  and the red highlighted region indicates the error bound described in Theorem 3. (a) and (b) show how each error metric changes over each training epoch for the models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ . Note that  $\mathcal{E}_{\text{approx}}$  is 0.00008 for the  $L_1$  regularized model and  $\mathcal{E}_{\text{approx}}$  is 0.212 for the un-regularized model. Additionally, (b) and (c) show how the model parameters update from epoch for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ , respectively.



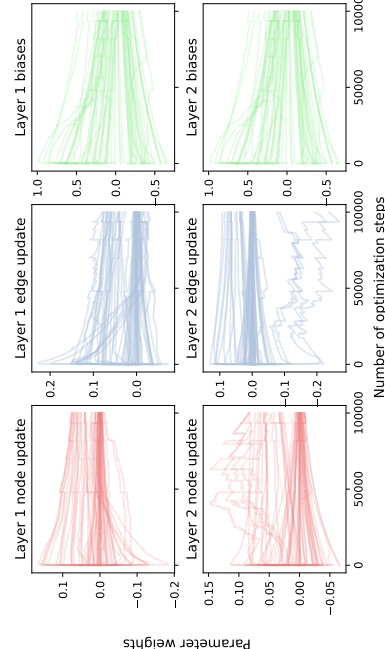
(a) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .



(b) Model parameters for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .



(c) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}}$ .



(d) Model parameters for models trained with  $\mathcal{L}_{\text{MSE}}$ .

Figure 9: Error metrics and parameter updates for a two-layer MinAgg GNN trained on a single step of the Bellman-Ford algorithm. Note that the dotted line in (a) and (b) is the global minimum of Eq. (8). (a) and (b) show how the train loss and test loss change over each training epoch for the models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$  as well as the theoretical loss term  $\mathcal{L}_{\text{reg}}$  over time.  $\mathcal{L}_{\text{approx}}$  converges to 0.001 for the  $L_1$  regularized model and 0.0312 for the un-regularized model. (b) and (c) show how the model parameters update from epoch for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ , respectively. Each curve has been smoothed with a truncated Gaussian filter with  $\sigma = 20$ .

## D.2. Deep MinAgg GNNs

We examine a variety of MinAgg GNN configurations to empirically verify Theorem 3. As in Section 3, we compare models trained with  $L_1$  regularization against models trained using just  $\mathcal{L}_{\text{MSE}}$ .

Similar to Section 3, for each model, we track  $\mathcal{E}_{\text{approx}}$ ,  $\mathcal{L}_{\text{train}}$ , and  $\mathcal{L}_{\text{reg}}$  throughout optimization. Note that  $\mathcal{E}_{\text{approx}}$  is evaluated on the same set of test graphs described in Section 3 i.e.  $\mathcal{G}_{\text{test}}$  is a set of 200 graphs which are a mix of cycles, complete graphs, and Erdős-Renyí graphs with  $p = 0.5$ . Note that Theorem 3 requires  $\mathcal{L}_{\text{reg}}$  to fall below a certain  $\epsilon$  threshold (indicated by the red region). For each of the complex GNN configurations we analyze below, we observe that models trained with  $L_1$  regularization satisfy this bound, aligning with Theorem 3. Furthermore, to verify that the model learns the correct parameters which implement Bellman-Ford, we also track a summary of the model parameters per epoch, defined as follows. This is the same as the model parameters from Section 3 in the main text, but here we provide more detail. Recall the definition of MinAgg GNNs in Theorem 1. Each layer can be precisely expressed as follows:

$$\sigma\left(W^{\text{up},(\ell)}\left(\min\{\sigma(W^{\text{agg},(\ell)}(h_u^{(\ell-1)} \oplus x_{(u,v)})) + b^{\text{agg},(\ell)}\} \oplus h_v^{(\ell-1)}\right)\right) + b^{\text{up},(\ell)}$$

To analyze parameter dynamics, we visualize the following for each optimization step and each layer:

- Layer  $\ell$  node parameters: Given a node feature in  $\mathbb{R}^d$ , the first  $d$  columns of  $W^{\text{agg},(\ell)}$ ,  $W^{\text{agg},(\ell)}[:, :d]$  scale the incoming neighboring node features. The contribution of incoming neighboring node features to the layer-wise output node feature for  $v$  can be summarized as follows:

$$\bigoplus_{j=1}^d \left( \bigoplus_{i=1}^{d_{\text{up}}} W^{\text{up},(\ell)}[i, :d_{\text{agg}}] \odot W^{\text{agg},(\ell)}[:, j] \right) \oplus \left( \bigoplus_{i=1}^{d_{\text{agg}}+d_{\text{up}}} W^{\text{up},(\ell)}[i, d_{\text{agg}} : d_{\text{agg}} + d_{\text{up}}] \right)$$

where  $\odot$  is the element-wise product and  $\oplus$  denotes concatenation.

- Layer  $\ell$  edge parameters: Similar to above, we know that the last column,  $d + 1$ , of  $W^{\text{agg},(\ell)}$  scales the incoming edge features. Therefore, the contribution of neighboring edges to the layer-wise output node feature can be summarized as

$$\bigoplus_{i=1}^{d_{\text{up}}} W^{\text{up},(\ell)}[i, :d_{\text{agg}}] \odot W^{\text{agg},(\ell)}[:, d + 1].$$

- Layer  $\ell$  biases: We track the bias terms for each layer as  $b^{\text{agg},(\ell)} \oplus b^{\text{up}}$ . Note that for the sparse implementation of Bellman-Ford that we describe previously, we require that the bias terms all converge to zero.

Finally, note that we summarize each model’s size generalization ability in Table 2 (analogous to Table 1 in the main text) in both the setting where we make a single forward pass through the model and when we use the model as a module and make repeated forward passes through the model. In Table 1, we see that each of the models trained with  $L_1$  regularization achieve significantly lower test error than those trained without  $L_1$  regularization.

### D.2.1. ONE LAYER

We configure the single layer GNN with BF-GNN with  $d_{\text{agg}} = 64$  and  $d_{\text{up}} = 1$ . We evaluate  $\mathcal{L}_{\text{MSE}}$  on  $\mathcal{G}_{\text{train}}$  which is the same as for the simple BF-GNN (four two-node path graphs and four three-node path graphs starting from step one of Bellman-Ford). Intuitively, the two-node path graphs provide a signal which controls the edge update feature while the three-node path controls the node update feature of Bellman-Ford. The results for a single layer BF-GNN are summarized in Fig. 8.

First, as illustrated in Fig. 8 (a) and (b), the train error  $\mathcal{L}_{\text{MSE}}$  alone does not capture the model’s generalization ability. Both the  $L_1$ -regularized and non-regularized models achieve low  $\mathcal{L}_{\text{MSE}}$  (0.002 for both the  $L_1$  regularized model and the un-regularized model), yet only the regularized model—where  $\mathcal{L}_{\text{reg}}$  converges to its minimum value—exhibits low test error.

Furthermore, we verify in Fig. 8 (b) and (c) that the parameters for the single layer MinAgg GNN converge to a configuration which *approximately implements* a single step of Bellman-Ford. Since we are only considering a single layer MinAgg GNN and the input node feature are the initial distances to source ( $d = 1$ ), the node parameter summary that we consider is

$$W^{\text{up},(1)}[: d_{\text{agg}}] \odot W^{\text{agg},(1)}[:, 1] \oplus W^{\text{up},(1)}[d_{\text{agg}} + d]$$

where  $W^{\text{agg},(1)} \in \mathbb{R}^{d_{\text{agg}} \times 2}$  and  $W^{\text{up},(1)} \in \mathbb{R}^{d_{\text{agg}} + d}$ . Additionally, the edge parameter summary we consider is

$$W^{\text{up},(1)}[: d_{\text{agg}}] \odot W^{\text{agg},(1)}[:, 2].$$

Therefore, for a single layer MinAgg GNN, the model parameters which exactly implement Bellman-Ford has a single  $k \in [64]$  such that

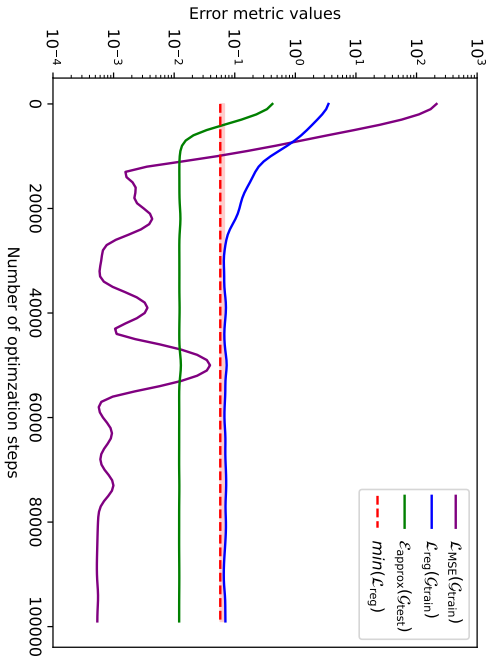
$$W^{\text{up},(1)}[k] \cdot W^{\text{agg},(1)}[k, 1] = W^{\text{up},(1)}[k] \cdot W^{\text{agg},(1)}[k, 2] = 1.0.$$

This means there is a single identical positive non-zero value for both the node and edge parameters. Additionally, all biases are zero. In Fig. 8 (c), we see that the trained MinAgg GNN using  $L_1$  regularization approximately converges to this configuration of parameters. However, the MinAgg GNN trained without regularization does not achieve this parameter configuration, explaining the higher test error  $\mathcal{E}_{\text{approx}}$ .

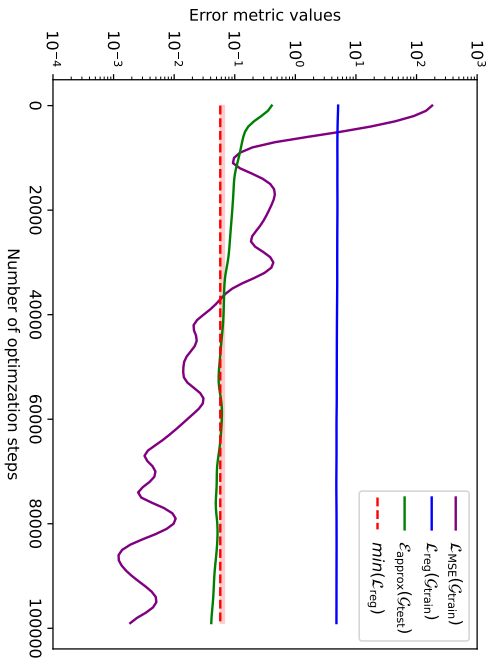
### D.2.2. TWO LAYER, SINGLE STEP

We configure a two layer GNN with  $d_{\text{agg}} = 64$  and  $d_{\text{up}} = 1$  for all layers and train it on a single step of Bellman-Ford. Note that this setup is overparameterized for modeling a single step of Bellman-Ford. The results of training on a single step of Bellman-Ford are summarized in Fig. 9. As with the single layer MinAggGNN configuration, the train set again consists of four two-node path graphs and four three-node path graphs with varying edge weights.

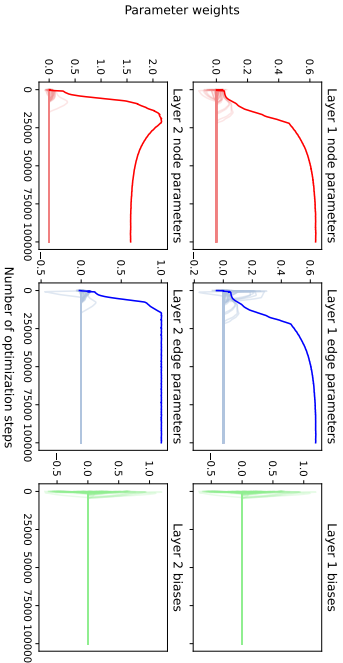
First, while  $\mathcal{L}_{\text{MSE}}$  converged to a low error for both models, only the model trained with  $\mathcal{L}_{\text{MSE}, L_1}$  has  $\mathcal{L}_{\text{reg}}$  converging to a low value. This also corresponds to a significantly lower  $\mathcal{L}_{\text{test}}$ , again demonstrating an overparameterized model’s ability to learn Bellman-Ford under an  $L_1$ -regularized training error and generalize to larger graph sizes with lower test error. These results align with our theoretical results, where the convergence of  $\mathcal{L}_{\text{reg}}$  to its minimum value provides a certificate for size generalization.



(a) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .



(b) Error metrics for models trained with  $\mathcal{L}_{\text{MSE}}$ .



(c) Model parameters summaries for model trained (d) Model parameters summaries for model trained with  $\mathcal{L}_{\text{MSE}, L_1}$ .

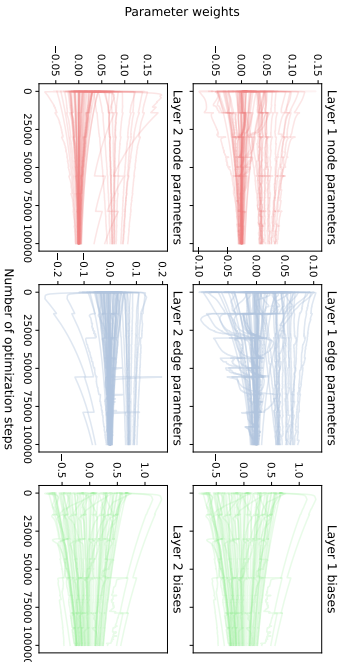


Figure 10: Performance metrics and parameter updates for a two-layer MinAgg GNN trained on two steps of the Bellman-Ford algorithm. The dotted line in (a) and (b) is the global minimum of Eq. (8) and the red region represents the  $\epsilon$  bounds in Theorem 3. In (a) and (b), we track the change in the train loss, test loss, and  $\mathcal{L}_{\text{reg}}$  over each optimization step for the models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ . Note that the final test loss for the model trained with  $\mathcal{L}_{\text{MSE}, L_1}$  is 0.006 while the final test loss for the model trained with  $\mathcal{L}_{\text{MSE}}$  is 0.288. (b) and (c) show how the model parameters change over each optimization step models trained with  $\mathcal{L}_{\text{MSE}, L_1}$  and  $\mathcal{L}_{\text{MSE}}$ , respectively. Each curve has been smoothed with a truncated Gaussian filter with  $\sigma = 20$ .

Table 2: Error ( $\mathcal{E}_{\text{approx}}$ ) versus size for all model configurations. The first row of the table contains the test error for both the single (indicated by 1L) and two-layer (indicated by 2L) model configurations trained on a single step of Bellman-Ford and the second row contains the test error for all model configurations trained on two steps of Bellman-Ford. We use ‘reg’ to indicate that the model is trained with  $L_1$  regularization and ‘un-reg’ to indicate a model trained without  $L_1$  regularization. Similar to Table 1, we examine the error for a single pass of each model (one step of BF for the first row and two steps of BF for the second row) and for three forward passes of each model (three steps of BF for the first row and six steps of BF for the second row). Each test set consists of Erdős–Rényi graphs generated with the corresponding sizes listed with  $p$  such that the expected degree  $np = 5$ . For both models trained with  $L_1$  regularization and without, the error does not change much as the number of nodes in the test graphs increase. However, when each model is used as a module and iterated, we see that the  $L_1$  regularized model remains accurate while the error for the un-regularized model increases significantly.

# steps for training	# of nodes	Single forward pass			
		1L, un-reg.	1L, reg.	2L, un-reg.	2L, reg.
One step	100	$0.0079 \pm 0.0027$	$0.00006 \pm 0.00003$	$0.0569 \pm 0.0064$	$0.0022 \pm 0.0002$
	500	$0.0070 \pm 0.0025$	$0.00006 \pm 0.00004$	$0.0560 \pm 0.0072$	$0.0022 \pm 0.0003$
	1K	$0.0071 \pm 0.0026$	$0.00006 \pm 0.00004$	$0.0558 \pm 0.0073$	$0.0021 \pm 0.0003$
Two steps	100	-	-	$0.0296 \pm 0.002$	$0.0173 \pm 0.001$
	500	-	-	$0.0297 \pm 0.002$	$0.0174 \pm 0.001$
	1K	-	-	$0.0308 \pm 0.002$	$0.0180 \pm 0.001$

# steps for training	# of nodes	Iterated, three forward passes			
		1L, un-reg.	1L, reg.	2L, un-reg.	2L, reg.
One step	100	$0.0320 \pm 0.0028$	$0.00012 \pm 0.00004$	$0.0097 \pm 0.0017$	$0.00133 \pm 0.0001$
	500	$0.0289 \pm 0.0025$	$0.00012 \pm 0.00006$	$0.0074 \pm 0.0010$	$0.00147 \pm 0.0001$
	1K	$0.0290 \pm 0.0025$	$0.00011 \pm 0.00001$	$0.0072 \pm 0.0013$	$0.00151 \pm 0.0001$
Two steps	100	-	-	$0.0596 \pm 0.0131$	$0.0182 \pm 0.0009$
	500	-	-	$0.0391 \pm 0.0040$	$0.0197 \pm 0.0006$
	1K	-	-	$0.0367 \pm 0.0025$	$0.0199 \pm 0.0007$

In Fig. 9 (b) and (c), we again further verify the size generalization ability of our models and emphasize the importance of sparsity regularization for size generalization by showing that the model parameter summaries for the model trained with  $\mathcal{L}_{\text{MSE}, L_1}$  converge to an implementation of a single step of Bellman-Ford. Since  $d_{\text{up},(\ell)} = 1$  for both layers, the node parameter summary that we analyze is again

$$W^{\text{up},(\ell)}[: d_{\text{agg}}] \odot W^{\text{agg},(\ell)}[:, 1] \oplus W^{\text{up},(\ell)}[d_{\text{agg}} + 1]$$

and the edge parameter summary is

$$W^{\text{up},(\ell)}[: d_{\text{agg}}] \odot W^{\text{agg},(\ell)}[:, 2].$$

for both layers. Similar to the single layer and single edge case, for the first layer, we expect that in the sparse implementation of a single step of Bellman, there will only a unique identical and positive non-zero value for both the node and edge parameter summaries. Therefore, for

this sparse implementation, for any node  $v$  in a given input graph, the node feature for  $v$  at the first layer is  $a(x_{u'} + x_{(u',v)})$  where  $a > 0$  and  $u' = \operatorname{argmin}_{u \in \mathcal{N}(v)} \{x_u + x_{(u,v)}\}$ . In the second layer,  $W^{\text{up},(2)}[65] = 1/a$  is the only positive non-zero parameter. Therefore, the final output for  $v$  will be  $\min_{u \in \mathcal{N}(v)} \{x_u + x_{(u,v)}\}$ . In Fig. 9 (b) and (c), we see again that the model trained with  $L_1$  regularization (i.e.  $\mathcal{L}_{\text{MSE}, L_1}$ ) has its parameters approximately converge to this sparse implementation of Bellman-Ford. In contrast, the model trained without  $L_1$  regularization (i.e. only  $\mathcal{L}_{\text{MSE}}$ ) does not appear to converge such a sparse implementation of Bellman-Ford, which accounts for the higher test error of the model.

### D.2.3. TWO LAYER, TWO STEPS

In the main text, we evaluate the ability of a two layer MinAgg GNN to learn two steps of Bellman-Ford. Here, we show the ability of the MinAgg GNN to learn two steps of Bellman-Ford in a somewhat under-parameterized setting as we let  $d_{\text{agg}} = 64$  and  $d_{\text{up}} = 1$ . The results are summarized in Fig. 10. Similar to the other model configurations evaluated (both in the supplement and the main text), we see that in Fig. 10 (a) and (b) that the  $L_1$  regularized model achieves much lower  $\mathcal{L}_{\text{reg}}$  and correspondingly, much lower  $\mathcal{E}_{\text{approx}}$ . The parameter configurations visualized in Fig. 10 (c) and (d) show that the  $L_1$  regularized model with low  $\mathcal{L}_{\text{reg}}$  converges to the sparse implementation of two-steps Bellman-Ford, as we have shown theoretically in Theorem 3. Additionally, note in Table 2, the gap between the error for the unregularized model and the error for the regularized model is much lower than that of Table 1 in the main text.

### D.3. Additional synthetic and real datasets

To provide further evidence of our claims, we provide a comparison of both the unregularized and  $L_1$ -regularized BF GNN error on several synthetic and real datasets. For the synthetic datasets we use:

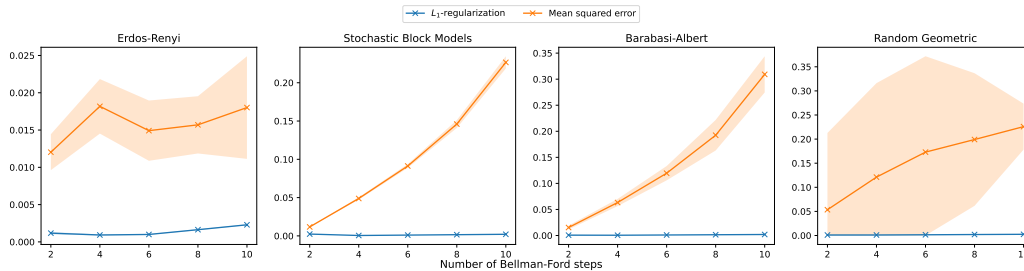
- Stochastic block models (SBM): these graphs are generated using 3-7 partitions of 35 vertices each and a randomly generated probability matrix.
- Barabasi-Albert models (BA): these graphs are generated according to the Barabasi-Albert model with preferential attachment parameter  $m \in [3, 65]$ .
- Random geometric graphs: we randomly position 5-65 nodes in  $[0, 1]^2$ . Two nodes are joined by an edge if the distance between them is less than 0.1.

For the real datasets, we use a terrain triangulation of patch of land in Norway with 40,000 vertices. Edge weights are determined by the Euclidean distance between points on the terrain. Additionally, we use Airports USA (USAir97) Rossi and Ahmed (2015); Rossi and Ahmed, where nodes represent airports and edges represent the existence of commercial flights between them. All edge weights for this graph are randomly generated. In Table 3, we evaluate the performance of the two-step BF GNN model. For the sake of comparison, we also include a comparison to the GAT (both trained with mean squared error and also trained with  $L_1$  regularization).

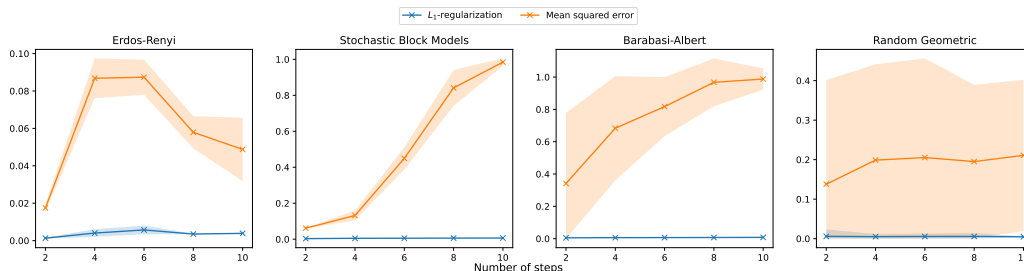
Additionally, we further evaluate our model as an iterative module for solving more than  $k$  steps of Bellman-Ford on several synthetic test graphs. We show this for both learning a single

	Synthetic			Norway	Airports-USA
	SBM	BA	Geometric		
MSE	0.053 ± 0.006	0.441 ± 0.431	0.226 ± 0.334	0.0141 ± 0.0652	0.1119 ± 0.0924
$L_1$ -reg	0.003 ± 0.0005	0.005 ± 0.0011	0.003 ± 0.0022	0.0003 ± 0.0008	0.0081 ± 0.0299
GAT (MSE)	0.977 ± 0.015	0.838 ± 0.157	0.971 ± 0.047	1.611 ± 3.027	5.426 ± 7.044
GAT ( $L_1$ )	0.722 ± 0.121	1.237 ± 0.2137	0.616 ± 0.297	0.9178 ± 0.1416	0.998 ± 0.0081

Table 3: Comparison of test error across all synthetic and real datasets.



(a) Iterating one step of Bellman-Ford



(b) Iterating two steps of Bellman-Ford

Figure 11: Extended results for iterative modules. We compare the model trained with an  $L_1$ -regularized loss to a model trained with mean squared error loss

step of Bellman–Ford and learning two steps of Bellman–Ford; see Fig. 11. Across iterations, the  $L_1$ -regularized models exhibit substantially slower error accumulation (flatter curves) than the unregularized  $\mathcal{L}_{\text{MSE}}$  models.

#### D.4. Hidden edge weights

We next evaluate a setting where the algorithmic skeleton (shortest paths via Bellman–Ford) is appropriate, but a component of the computation is unknown and must be learned. Concretely, edges expose an observed attribute  $x_e$ . The effective distance used by the shortest-path objective is an unknown, nonnegative transformation  $g(x_e)$ . The model must both learn  $g$  and learn to execute BF, the latter of which is aided by algorithmic alignment. As a simple instantiation, we set  $g(x) = (x - 2)^2$  and supervise on shortest-path labels computed with these transformed edge costs (see Table 4).

This scenario mirrors real applications—for example, routing to minimize fuel or energy consumption where per-road cost is a nonlinear function of observable features (e.g., grade, speed limit, and congestion) rather than raw length alone. In this regime, the  $L_1$ -regularized, BF-aligned MinAgg GNN reliably learns a sparse approximation to  $g$  and extrapolates to larger and structurally different graphs, substantially outperforming the unregularized  $\mathcal{L}_{\text{MSE}}$  model. This provides a concrete example of how algorithmic alignment remains effective when we have only a prior over the desired computation (e.g., shortest-path computation) and must infer the remaining components from data.

	ER	SBM	BA	Geometric
MSE	$0.383 \pm 0.055$	$0.541 \pm 0.093$	$0.276 \pm 0.129$	$0.517 \pm 0.732$
$L_1$ -reg	$0.0054 \pm 0.0006$	$0.0055 \pm 0.0006$	$0.0057 \pm 0.0013$	$0.0291 \pm 0.0334$

Table 4: Hidden edge function  $f(x) = (x - 2)^2$

## Appendix E. Limitations

We now discuss limitations to our current study. First, we do not provide experiments for graph neural networks (GNNs) with many layers, and thus our experimental findings may not generalize to deeper architectures. Second, our approach requires training sets to be explicitly constructed rather than sampled from a distribution, placing it closer in spirit to active query learning or machine teaching than to standard i.i.d. learning. Our goal, however, is to prove OOD generalization from such examples, not to optimize or minimize the teaching set. Finally, we focus exclusively on the properties of the global minimum of the loss and do not discuss optimization dynamics or the process of reaching this minimum, which are important considerations in real-world training.

## Appendix F. Table of Notation

Symbol	Definition
$[n]$	The set $\{1, 2, \dots, n\}$ .
$x \oplus y$	Concatenation of vectors $x$ and $y$ .
$x_i$ or $[x]_i$	$i$ -th component of vector $x$ .
$\beta$	A large constant representing the unreachable node feature value.
<b>Graphs and Attributed Graphs</b>	
$G = (V, E, X_v, X_e)$	Attributed graph with vertices $V$ , edges $E$ , edge weights $X_e$ , and node attributes $X_v$ .
$X_e, X_v$	Edge weights $\{x_e : e \in E\}$ and node attributes $\{x_v : v \in V\}$ .
$d^{(t)}(s, v)$	Length of the $t$ -step shortest path from node $s$ to $v$ ; $\beta$ if no such path exists.
$G^{(t)}$	$t$ -step Bellman-Ford (BF) instance with node features $\{x_v = d^{(t)}(s, v) : v \in V\}$ .
$\Gamma$	Operator implementing a single step of the BF algorithm.
$\mathcal{G}$	Set of all edge-weight-bounded attributed graphs.
$P_k^{(\ell)}(a_1, \dots, a_k)$	$k$ -edge path graph at step $\ell$ with edge weights $a_1, \dots, a_k$ .
$\mathcal{N}(v)$	Neighborhood of node $v$ in the graph.
$V^*(G)$	Set of reachable nodes in graph $G$ , i.e., nodes with $x_v \neq \beta$ .
<b>Graph Neural Networks (GNNs)</b>	
$\mathcal{A}_\theta$	$L$ -layer Bellman-Ford Graph Neural Network (MinAgg GNN) parameterized by $\theta$ .
$h_v^{(\ell)}$	Hidden feature of node $v$ at layer $\ell$ in the MinAgg GNN.
$\mathcal{A}_\theta(G)$	Output graph of MinAgg GNN $\mathcal{A}_\theta$ after $L$ layers, with updated node features.
$d_\ell$	Dimensionality of hidden features at layer $\ell$ .
$f^{\text{agg}}, f^{\text{up}}$	MLPs used for aggregation and update operations in MinAgg GNN layers.
$W_j^{\text{agg}}, W_j^{\text{up}}$	Weight matrices for aggregation and update MLPs in MinAgg GNN.
$b_j^{\text{agg}}, b_j^{\text{up}}$	Bias vectors for aggregation and update MLPs in MinAgg GNN.
$K$	Number of message passing steps in the MinAgg GNN.
$m$	Number of layers in the MLPs used for aggregation and update in MinAgg GNN.
$L$	Total number of layers in the MinAgg GNN.
$d$	Dimensionality of hidden features in MinAgg GNN layers.
<b>Training and Loss Functions</b>	
$\mathcal{H}_{\text{small}}$	A set of small training graphs used to analyze GNN performance.
$\mathcal{G}_{\text{train}}$	Set of training examples for the MinAgg GNN, consisting of input-output graph pairs.
$\mathcal{L}_{\text{reg}}$	Regularized loss function for MinAgg GNN, combining training loss and parameter sparsity penalty.
$\mathcal{L}_{\text{MAE}}$	Mean absolute error loss over the training set.
$\eta$	Regularization coefficient for sparsity in the MinAgg GNN.
$\mathcal{E}_{\text{approx}}$	Average relative error of model predictions vs. ground-truth BF output over the test set.
$\mathcal{L}_{\text{MSE}}$	Mean Squared Error.

**Appendix G. LLM Use**

LLMs were used for refining writing, proof-reading, and organizing citations.

**Appendix H. Table of Contents****Contents**

<b>1 Introduction</b>	<b>1</b>
1.1 Related work . . . . .	4
<b>2 Extrapolation Guarantees</b>	<b>5</b>
2.1 Model . . . . .	5
2.2 Warm up . . . . .	6
2.3 Main Result . . . . .	7
<b>3 Experiments</b>	<b>9</b>
<b>4 Discussion</b>	<b>12</b>
<b>A Definitions and notation</b>	<b>17</b>
<b>B Warm-up: Single layer GNNs implement one step of BF</b>	<b>19</b>
<b>C Sparsity regularized deep GNNs implement BF</b>	<b>26</b>
C.1 Implementing BF . . . . .	27
C.2 Training set . . . . .	28
C.3 Sparsity structure . . . . .	29
C.4 Bounding GNN expressivity . . . . .	35
C.5 Global minimum is BF . . . . .	39
<b>D Additional Experiments</b>	<b>46</b>
D.1 Simple MinAgg GNN . . . . .	46
D.2 Deep MinAgg GNNs . . . . .	50
D.2.1 One layer . . . . .	51
D.2.2 Two layer, single step . . . . .	51
D.2.3 Two layer, two steps . . . . .	54
D.3 Additional synthetic and real datasets . . . . .	54
D.4 Hidden edge weights . . . . .	55
<b>E Limitations</b>	<b>56</b>

<b>F</b>	<b>Table of Notation</b>	<b>57</b>
<b>G</b>	<b>LLM Use</b>	<b>58</b>
<b>H</b>	<b>Table of Contents</b>	<b>58</b>