

ADMP-GNN: Adaptive Depth Message Passing GNN

Yassine Abbahaddou
LIX, Ecole Polytechnique, IP Paris
Palaiseau, France
yassine.abbahaddou@polytechnique.edu

Johannes F. Lutzeyer
LIX, Ecole Polytechnique, IP Paris
Palaiseau, France
johannes.lutzeyer@polytechnique.edu

Fragkiskos D. Malliaros
Centralesupelec, Universite Paris-Saclay, Inria
Gif-sur-Yvette, France
fragkiskos.malliaros@centralesupelec.fr

Michalis Vazirgiannis
LIX, Ecole Polytechnique, IP Paris
Palaiseau, France
mvazirg@lix.polytechnique.fr

Abstract

Graph Neural Networks (GNNs) have proven to be highly effective in various graph learning tasks. A key characteristic of GNNs is their use of a fixed number of message-passing steps for all nodes in the graph, regardless of each node's diverse computational needs and characteristics. Through empirical real-world data analysis, we demonstrate that the optimal number of message-passing layers varies for nodes with different characteristics. This finding is further supported by experiments conducted on synthetic datasets. To address this, we propose *Adaptive Depth Message Passing GNN (ADMP-GNN)*, a novel framework that dynamically adjusts the number of message passing layers for each node, resulting in improved performance. This approach applies to any model that follows the message passing scheme. We evaluate ADMP-GNN on the node classification task and observe performance improvements over baseline GNN models. Our code is publicly available at: <https://github.com/abbahaddou/ADMP-GNN>.

CCS Concepts

• **Theory of computation** → **Randomness, geometry and discrete structures**; **Randomness, geometry and discrete structures**; • **Computing methodologies** → **Machine learning**.

Keywords

Graph Neural Networks, Adaptive Neural Networks

ACM Reference Format:

Yassine Abbahaddou, Fragkiskos D. Malliaros, Johannes F. Lutzeyer, and Michalis Vazirgiannis. 2025. ADMP-GNN: Adaptive Depth Message Passing GNN. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management (CIKM '25)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '25, Coex, Seoul, Korea

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

A plethora of structured data comes in the form of graphs [3, 5]. This has driven the need to develop neural network models, known as Graph Neural Networks (GNNs), that can effectively process and analyze graph-structured data. GNNs have garnered significant attention for their ability to learn complex node and graph representations, achieving remarkable success in several practical applications [6, 9, 10, 31]. Many of these models are instances of Message Passing Neural Networks (MPNNs) [16, 24, 40]. A common characteristic of GNNs is that they typically employ a fixed number of message passing steps for all nodes, determined by the number of layers in the GNN [16]. This static framework raises an intriguing question: *Should the number of message passing steps be adapted individually for each node to better capture their unique characteristics and computational needs?*

Determining the optimal number of message passing layers for each node in a GNN presents a significant challenge due to the intricate and diverse nature of graph structures, node features, and learning tasks. While deeper GNNs can capture long-range dependencies [28], they can also encounter issues like oversmoothing, where nodes become indistinguishably similar [17, 29]. In dense graphs, where information can propagate quickly, even shallow GNNs can effectively capture local information [43]. Conversely, sparse graphs, particularly those with isolated nodes or limited connectivity, may require additional layers to facilitate effective information sharing [44, 45]. This underscores the importance of selecting the appropriate number of layers for a GNN to capture the necessary graph information effectively. An even more compelling idea is to adjust the GNN depth for each node based on its local structural properties. This adaptive approach could be especially beneficial for graphs with varied local structures, ensuring that each node is processed according to its unique requirements.

This need for per-node customization naturally aligns with the concept of *Dynamic Neural Networks*, also referred to as *Adaptive Neural Networks* [2]. Dynamic Neural Networks represent a class of models that can adjust their architecture or parameters depending on the input. Dynamic Neural Networks have gained significant popularity, especially in the field of computer vision. Examples of adaptation include varying the number of layers and implementing skip connections [20, 25, 33]. Beyond computer vision, other types of adaptive neural networks have been explored in various domains. In natural language processing, adaptive computation time models allow recurrent neural networks (RNNs) to determine the required

number of recurrent steps dynamically based on the complexity of the input sequence [18]. However, applying these adaptations to graph learning tasks presents unique challenges. Unlike the structured and homogeneous data often encountered in computer vision, graph data involves overcoming nuisances related to the inherent complex structure of graphs. Graphs may have varying node degrees, non-uniform connectivity, and feature Heterogeneity. While some dynamic approaches may extend effectively to graph classification tasks, applying these methods to node classification presents additional complexities. In node classification, each input sample (node) is part of a larger interconnected system where information propagates through edges, i.e., the prediction for one node often depends on the features of the other nodes and the structure of the graph. As a result, specific techniques are needed to account for the dependencies and relational information encoded in the graph structure.

In this work, we focus on the task of node classification by proposing ADMP-GNN, a novel approach that dynamically adapts the number of layers for each node within a GNN. Our main contributions are as follows:

- **Node-specific depth analysis in GNNs.** We demonstrate through empirical analysis that different nodes within the same graph may require varying numbers of message passing steps to accurately predict their labels. This finding underscores the importance of node-specific depth in GNNs.
- **Adaptive message passing layer integration.** We present ADMP-GNN, a novel approach that enables any GNN to make predictions for each node at every layer. Training the GNN to predict labels across all layers is a multi-task setting, which often suffers from gradient conflicts, leading to suboptimal performance. To address this, we propose a sequential training methodology where layers are progressively trained, and their gradients are subsequently frozen, thereby mitigating conflicts and improving overall performance.
- **Adaptive layer policy learning for node classification.** We introduce a heuristic method to learn a layer selection policy using a set of validation nodes. This policy is then applied to select the optimal layer for predicting the labels of test nodes, ensuring that each node exits the GNN at the most appropriate layer for its specific classification task.
- **Model-agnostic flexibility.** Our approach is model-agnostic and can be integrated with any GNN architecture that employs a message passing scheme. This flexibility enhances the GNN's performance on node classification tasks, providing a significant improvement over traditional fixed-layer approaches.

2 Related Work

In this section, we review key developments in GNNs and dynamical neural networks, which form the foundation of our work.

2.1 Graph Neural Networks

Graphs, denoted as $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$, are mathematical structures used to represent pairwise relationships. Here, \mathcal{V} is the set of nodes, \mathcal{E} is the set of edges, and $X = \{x_u : u \in \mathcal{V}\}$ represents the node features, where each $x_u \in \mathbb{R}^d$ corresponds to a d -dimensional

feature vector associated with node u . GNNs are a class of deep learning methods specifically designed to operate on such graph-structured data. As with most neural networks, GNNs are formed by stacking many layers. Each layer, ℓ , is responsible for updating the node representations $\{h_u^{(\ell)} : u \in \mathcal{V}\}_{0 \leq \ell \leq L}$, relying on the graph structure and the output from the previous layer, $\{h_u^{(\ell-1)} : u \in \mathcal{V}\}$. Conventionally, the features of the nodes are used as input of the first layer, i.e., $h_u^{(0)} = x_u \in \mathbb{R}^d$ for all $u \in \mathcal{V}$. A basic GNN layer is based on the message passing mechanism and consists of two components: (i) *Aggregate Layer* ψ that applies for each node v , a permutation invariant function to its neighbors, denoted by $\mathcal{N}(v)$ to generate the aggregated node feature; (ii) *Update Layer* ϕ that combines the aggregated node feature $m_v^{(\ell)}$ with the previous hidden vector $h_v^{(\ell-1)}$, and generate a new representation $h_v^{(\ell)}$ of the same node v ,

$$\begin{aligned} m_v^{(\ell)} &= \psi^{(\ell)}(\{h_u^{(\ell-1)} : u \in \mathcal{N}(v)\}), \\ h_v^{(\ell)} &= \phi^{(\ell)}(h_v^{(\ell-1)}, m_v^{(\ell)}). \end{aligned}$$

Depending on the task, an additional readout or pooling function can be added after the last layer to aggregate the representation of nodes, as follows, $h_G = \text{ReadOut}(\{h_u^{(L)} : u \in \mathcal{V}\})$. Examples of the most popular GNN models include the Graph Convolutional Network (GCN) [24] and Graph Isomorphism Network (GIN) [40], which apply a weighted sum and sum operator in their aggregate layer, respectively.

2.2 Dynamical Neural Networks

Dynamic neural networks are gaining significance in the field of deep learning. Unlike static models with fixed computational graphs and parameters during inference, dynamic networks adapt their structures or parameters based on varying inputs. This dynamic flexibility gives models significant benefits, such as improved accuracy, enhanced computational efficiency, and superior adaptability [39, 46]. A popular type of dynamic neural networks includes those that dynamically adjust network depth based on each input. For instance, in natural language processing, some adaptive large language models employ adaptive depth to optimize both inference speed and computational memory usage of the Transformer architecture [11, 34, 38]. In the field of computer vision, there are studies that dynamically generate filters conditioned on each input, enhancing flexibility without significantly increasing the number of model parameters [22].

In Graph Machine Learning, dynamic adaptations have been proposed for GNN message passing. These methods include dynamically determining which neighbors to consider at each layer, enabling more flexible and adaptive message passing [15], or allowing nodes to react to individual messages at varying times rather than processing aggregated neighborhood information synchronously [14]. Another line of work focuses on adapting normalization layers for each node to enhance expressive power, generating representations that reflect local neighborhood structures [12]. Additionally, other related approaches use residual connections to mitigate issues like oversmoothing [7, 13]. To the best of our knowledge, in the field of GNNs, there has been no prior work proposing adaptive depth for each node. However, several studies have focused on combining

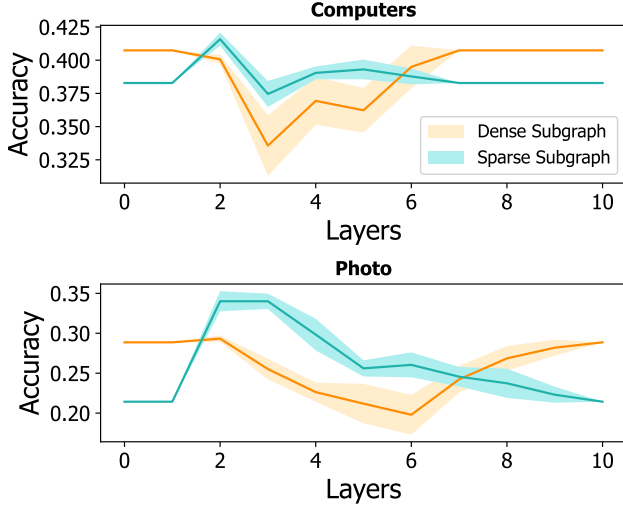


Figure 1: Effect of GCN’s depth on sparse and dense subgraphs. The figure shows the performance of GCNs when varying layer depths, and comparing its effectiveness on both sparse and dense subgraphs.

all GNN layers. These works typically aim to adapt GNN architectures for heterophilic graphs [8] and leverage information from higher-order neighbors [41]. While combining GNN layers can be viewed as a form of depth-adaptive strategy, where the final node representation is guided by the optimal intermediate hidden states, this approach remains *static* because the same inference policy is applied uniformly across all nodes and learned layer aggregators stay fixed after training.

3 Adaptive Depth Message Passing GNN

In this section, we first present an empirical analysis highlighting the necessity for node-specific depth in GNNs. Then, we introduce our ADMP-GNN. Our study includes experiments on both synthetic and real-world datasets to illustrate the importance and potential benefits of this methodology.

3.1 Depth Analysis on Synthetic Graphs

This analysis aims to highlight the importance of employing a varying number of message passing steps based on the specific characteristics of individual nodes. As outlined in the introduction, this approach is particularly relevant for graphs where nodes exhibit diverse properties, such as varying local structures. In this experiment, we investigate the effect of the number of message-passing layers on nodes with varied levels of local neighborhood sparsity. To achieve this, we construct a graph by merging two subgraphs extracted from real-world datasets, namely Computers and Photo [36]. Both subgraphs contain the same number of nodes, exhibit nearly identical homophily, and have equally distributed node labels. The main difference between these subgraphs lies in their structure, as one subgraph is sparse while the other is dense. Consequently, nodes within each subgraph share comparable structural characteristics. Details on the construction of these synthetic

datasets and visualizations of their adjacency matrices are provided in Appendix B.

We have trained $L + 1$ different GCN models, with a varying number of layers $\ell \in \{0, \dots, L\}$, where L represents the maximum depth, set to $L = 10$. Each GCN was trained on the entire synthetic graph, composed of both sparse and dense subgraphs, but the performance was evaluated separately on the individual subgraphs. This allows us to assess the impact of GNN depth on different types of local subgraph structures. The results, presented in Fig. 1, reveal notable differences in behavior between the subgraph types. As observed, in dense subgraphs, the test accuracy decreases at a faster rate, while in sparse subgraphs, the drop in accuracy occurs later, typically around layers 2 or 3. Moreover, the optimal number of layers differs between sparse and dense subgraphs. For instance, in the Computers dataset, the highest accuracy is achieved at layer 2 for the sparse subgraph, while for the dense subgraph, the optimal performance is reached at layer 0. Additionally, in the Photo dataset, we observe a distinct behavior starting from layer 6, where the impact of GNN depth diverges between sparse and dense subgraphs. This highlights the need to adapt the number of layers per node based on its characteristics.

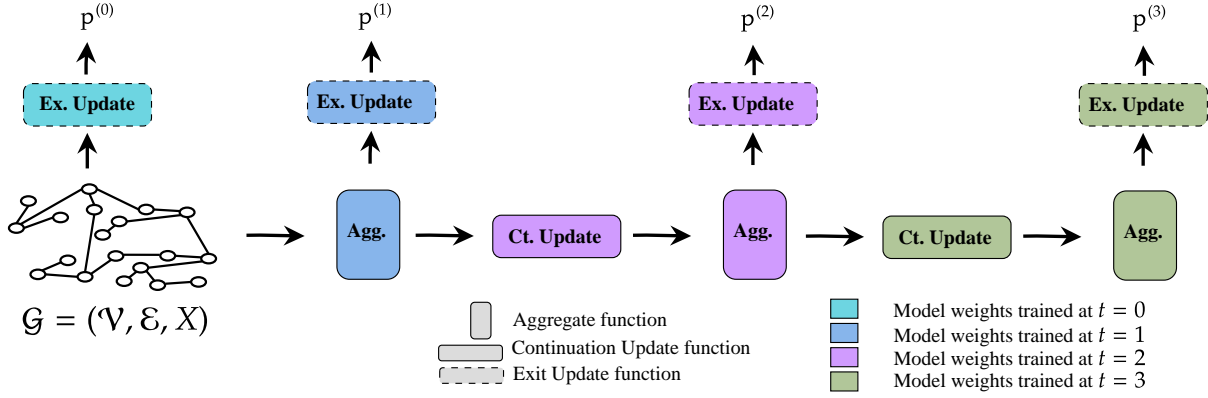
3.2 Adaptive Message Passing Layer Integration

Building on the insights from the previous analysis, it is essential to develop a framework that can efficiently adjust the GNN depth. For a maximum GNN depth L , a traditional approach would involve training $L + 1$ different GNNs, each with a distinct number of layers ℓ , where ℓ ranges from 0 and L . Subsequently, a policy must be established to determine the optimal GNN with the appropriate number of layers for each individual node. However, training $L + 1$ GNNs separately can be computationally expensive. A more efficient approach involves designing a single GNN with $L + 1$ layers that provide predictions at each intermediate layer. To ensure equivalence to the previous approach, i.e., training $L + 1$ GNNs separately, the computational graph for predictions at layer ℓ must match that of a standard GNN with ℓ layers, and the classification performance at layer ℓ should yield results comparable to those of a conventional GNN with ℓ layers.

To address these challenges, we introduce ADMP-GNN, an adaptation of a Message Passing Neural Network with a maximum depth of L layers. The goal is to ensure that the computational graph and the performance of ADMP-GNN at a certain layer match that of traditional GNNs when trained and tested on the same number of layers. To achieve this, we incorporate an additional *Update* function, denoted as $\phi_{\text{Ex}}^{(\ell)}$ (Ex stands for ‘Exit’), to directly predict node labels at a given layer ℓ , i.e., $p^{(\ell)}$. The function $\phi_{\text{Ex}}^{(\ell)}$ is defined as follows,

$$p_v^{(\ell)} = \phi_{\text{Ex}}^{(\ell)} \left(h_v^{(\ell-1)}, m_v^{(\ell)} \right) = \text{Softmax} \left(\tilde{W}^{(\ell)} m_v^{(\ell)} \right),$$

where $\tilde{W}^{(\ell)} \in \mathbb{R}^{d^{(\ell)} \times c}$ is a learnable weight matrix, $d^{(\ell)}$ is the dimension of the hidden representation at the ℓ -th layer, and c is the number of classes. To obtain predictions at a deeper layer $\ell' \geq \ell$, we continue the message passing using another *Update* function

Figure 2: Illustration of ADMP-GNN, when the maximum GNN depth is $L = 3$.

$\phi_{\text{Ct}}^{(\ell)}$ (Ct stands for ‘Continuation’),

$$p_v^{(\ell)} = \phi_{\text{Ex}}^{(\ell)}(h_v^{(\ell-1)}, m_v^{(\ell)}),$$

$$h_v^{(\ell+1)} = \phi_{\text{Ct}}^{(\ell)}(h_v^{(\ell-1)}, m_v^{(\ell)}).$$

For $\ell = 0$, we directly use the *Exit Update* function on the node features, i.e., $m_v^{(0)} = x_v$,

$$\forall v \in \mathcal{V}, \quad p_v^{(0)} = \phi_{\text{Ex}}^{(0)}(m_v^{(0)}) = \text{Softmax}(\tilde{W}^{(0)} x_v).$$

In Fig. 2, we illustrate the architecture of the proposed ADMP-GNN.

3.3 Training Scheme of ADMP-GNN

Our next objective is to train ADMP-GNN to predict node labels across all layers $\ell \in \{0, \dots, L\}$ simultaneously. For each layer ℓ , let θ_ℓ denote the weights of the function $\psi^{(\ell)} \circ \phi_{\text{Ct}}^{(\ell)}(\cdot)$. Two distinct strategies for training ADMP-GNN are explored:

1. *Aggregate Loss Minimization (ALM)*. This straightforward approach minimizes the aggregate loss over all layers. The total loss is formulated as,

$$\mathcal{L}_{\text{ALM}} := \arg \min_{\theta} \mathbb{E}_{v \in \mathcal{V}} [S_L(v)] \quad (1)$$

$$= \arg \min_{\theta_0, \dots, \theta_L} \mathbb{E}_{v \in \mathcal{V}} \left[\sum_{\ell=0}^L \mathcal{L}(p_v^{(\ell)}(m_v^{(0)}, \theta_0, \dots, \theta_\ell), y_v) \right],$$

where $p_v^{(\ell)}(m_v^{(0)}, \theta_0, \dots, \theta_\ell) = \phi_{\text{Ex}}^{(\ell)}(m_v^{(\ell)})$ is the prediction for the node v at layer ℓ , and \mathcal{L} is the Cross Entropy Loss. This approach may encounter gradient conflicts, particularly for early layers involved in both computation and back-propagation across upper layers.

2. *Sequential Training (ST)*. We have studied an alternative training setup where we progressively train one GNN layer at a time, subsequently freezing each layer after training. More formally, the

problem in (1) can be tackled using dynamic programming as follows,

$$\forall v \in \mathcal{V}, S_{\ell+1}(v) = \mathcal{L}(p_v^{(\ell)}(m_v^{(0)}, \theta_0^*, \dots, \theta_\ell^*, \theta_{\ell+1}), y_v) + S_\ell(v)$$

$$\theta_\ell^* = \arg \min_{\theta_\ell} \mathbb{E}_{v \in \mathcal{V}} [S_\ell(v)],$$

where $\forall v \in \mathcal{V}$, $S_0(v) = \mathcal{L}(\phi_{\text{Ex}}^{(0)}, y_v)$. For each intermediate layer $\ell < L$, by training this layer on the node classification task, we obtain high-quality node representations $\{h_v^{(\ell)} : v \in \mathcal{V}\}$. These representations are directly employed for predictions and serve as a robust foundation for the label predictions of the subsequent layer $\ell + 1$. Algorithm 1 offers a summary of the approach.

Comparison of ADMP-GNN Training Paradigms. To identify the optimal multi-task training configuration, we evaluate how much of a performance drop we lose at each layer compared to the single-task setting in GNNs. The comparative analysis of the three strategies for both GCN is detailed in Tables 1. Our findings indicate that ADMP-GNN ST outperforms ADMP-GNN ALM. Notably, the performance of ADMP-GNN ST is comparable to, or even exceeds, that of GNN when trained under the single-task setting. Furthermore, ADMP-GNN ST exhibits a smaller standard deviation, suggesting more consistent performance. These results underscore the effectiveness of the ST setting in addressing gradient conflicts inherent in aggregate loss minimization (ALM). Furthermore, ADMP-GNN ST effectively mimics the results of training $L + 1$ separate GNNs, each with a different number of layers, while requiring only a single unified model. The next step, outlined in Section 3.5, is to learn a policy that selects the optimal prediction layer for each node, completing the framework of ADMP-GNN.

Time Complexity. The training setup ST, where we sequentially train the deep ADMP-GNN, incurs relatively higher time costs due to the need for $L + 1$ training iterations. However, in each iteration, backpropagation is performed on a limited number of parameters, approximately equivalent to those in a single message passing layer. Consequently, only a small number of epochs are required for each training iteration. We report the training time of each approach in Table 6 in Appendix C.

Table 1: Comparison of ADMP-GCN training paradigms ALM and ST. These paradigms are also compared to the single-task training setting to evaluate which approach most closely mimics the classical GCN under single-task training. The best results for each dataset are bolded.

# Layers	Training Paradigm	Model	Cora	CiteSeer	CS	PubMed	Genius	Ogbn-arxiv
①	Single-task	GCN	56.38±0.04	57.18±0.12	88.04±0.49	72.50±0.09	80.82±1.00	48.88±0.06
	Multi-task	ADMP-GCN (ALM)	56.96±0.20	58.44±0.21	87.06±1.06	72.11±0.18	80.03±0.37	36.50±0.12
		ADMP-GCN (ST)	56.38±0.06	57.17±0.09	87.27±1.29	72.48±0.14	80.17±0.79	48.86±0.03
①	Single-task	GCN	76.90±0.14	69.68±0.06	91.74±0.80	76.63±0.13	80.23±0.37	55.21±0.50
	Multi-task	ADMP-GCN (ALM)	75.67±0.18	70.12±0.04	90.55±0.74	73.74±0.16	80.13±0.29	39.54±1.44
		ADMP-GCN (ST)	76.90±0.00	69.70±0.00	90.89±0.81	76.60±0.00	79.93±0.00	55.15±0.00
②	Single-task	GCN	81.06±0.50	71.05±0.48	91.67±0.94	79.46±0.31	79.88±0.51	66.92±0.67
	Multi-task	ADMP-GCN (ALM)	70.82±4.05	60.95±2.39	31.20±12.85	75.10±2.41	79.64±0.60	55.25±0.92
		ADMP-GCN (ST)	80.73±0.33	71.33±0.40	91.49±0.66	79.02±0.21	80.06±0.11	66.51±0.65
③	Single-task	GCN	79.14±1.58	66.33±1.35	89.80±0.87	78.50±0.68	80.00±0.04	67.33±0.55
	Multi-task	ADMP-GCN (ALM)	68.64±5.14	51.30±6.74	47.82±12.98	74.17±1.98	80.04±0.10	56.22±0.50
		ADMP-GCN (ST)	80.21±0.52	70.08±0.90	89.83±1.02	78.25±0.51	79.93±0.00	68.31±0.48
④	Single-task	GCN	75.96±1.93	60.33±2.38	78.90±22.33	76.59±0.98	80.01±0.04	65.49±0.99
	Multi-task	ADMP-GCN (ALM)	67.88±6.05	49.29±7.87	52.76±11.79	73.40±1.95	80.04±0.10	56.43±0.31
		ADMP-GCN (ST)	81.05±0.49	67.99±0.74	88.86±0.70	75.27±1.02	79.93±0.00	69.29±0.74
⑤	Single-task	GCN	70.09±4.01	57.40±3.43	77.96±15.24	74.32±3.66	80.01±0.04	63.04±1.33
	Multi-task	ADMP-GCN (ALM)	67.68±7.04	50.14±7.65	54.53±10.61	73.27±1.51	80.04±0.10	56.66±0.31
		ADMP-GCN (ST)	81.22±0.36	67.42±0.75	86.65±1.52	75.08±0.94	79.93±0.00	68.92±1.00

3.4 Empirical Insights into Node Specific Depth

We define *Oracle Accuracy* as the maximum achievable test accuracy under an optimal policy for selecting layers. This is formally expressed as,

$$\mathcal{A}_{\text{oracle}} = \frac{1}{|\mathcal{V}_{\text{test}}|} \sum_{v \in \mathcal{V}_{\text{test}}} \mathbb{1} \left(y_v \in \{\hat{y}_v^{(\ell)} : \ell = 0, \dots, L\} \right),$$

where $\mathcal{V}_{\text{test}}$ is the set of test nodes, $\{\hat{y}_v^{(\ell)} : \ell = 0, \dots, L\}$ represents the predictions for node v at each layer, and $\mathbb{1}(\cdot)$ is the indicator function. Specifically, $\mathcal{A}_{\text{oracle}}$ corresponds to the accuracy obtained if, for each test node, we could perfectly choose the layer that provides the correct prediction. Assuming such an optimal selection for each node, it represents the upper bound of accuracy achievable when employing adaptive strategies to determine the best layer for each node.

Based on the findings presented in Table 2 for GCN, it is clear that $\mathcal{A}_{\text{oracle}}$ outperforms the highest accuracies achieved by both GCN and ADMP-GCN ALM. This empirical evidence strongly indicates that adopting a distinct exit layer for each node is highly beneficial in an optimal configuration.

3.5 Generalization to Test Nodes

In this section, we introduce a heuristic approach to predict the optimal exit layer for test nodes based on the assumption that nodes exhibiting *structural similarity* should share the same exit layer. The

notion of structural similarity can be assessed using various metrics. In our work, we define the structural similarity based on node centrality metrics, i.e., nodes are considered structurally similar if they exhibit closely aligned centrality values within the graph. We measured node centralities using both local metrics like degree and global metrics such as k -core [30], PageRank scores [4], and Walk Count, which are detailed below.

k -core. The k -core decomposition of a graph is a method that involves systematically pruning nodes. In each iteration, a subgraph $\mathcal{G}_k \subset \mathcal{G}$ is constructed by removing all vertices whose degree is less than k . The core number of a vertex u , denoted as $\text{core}(u)$, represents the largest value of k for which u is included in the k -core. Formally, this is expressed as $\text{core}(u) = \max\{k : u \in \mathcal{G}_k\}$. Nodes with higher core numbers are typically found in denser subgraphs, reflecting strong interconnections with their neighbors and indicating a higher level of centrality.

PageRank. A node’s PageRank score represents the probability of a random walk visiting that node, making it a fundamental metric for measuring node importance in social network analysis and the web [4].

Walk Count. The Walk Count centrality quantifies the importance of a node based on the number of walks of a given length ℓ starting from that node. This measure captures higher-order connectivity patterns within the graph, going beyond direct neighbors to account for paths that traverse intermediate nodes. In this work, we

Table 2: Comparison of highest accuracy (\pm standard deviation) for GCN and ADMP-GCN ST, with the layer achieving the best accuracy indicated in brackets. The final row shows the *Oracle Accuracy* for ADMP-GCN ST. Best results for each dataset are bolded.

Model	Cora	CiteSeer	CS	PubMed	Genius	Ogbn-arxiv
GCN	81.06 \pm 0.50 [2]	71.05 \pm 0.48 [2]	91.67 \pm 0.94 [2]	79.46 \pm 0.31 [2]	80.82 \pm 1.00 [0]	67.33 \pm 0.55 [3]
ADMP-GCN (ST)	81.22 \pm 0.36 [5]	71.33 \pm 0.40 [2]	91.49 \pm 0.66 [2]	79.02 \pm 0.21 [2]	80.17 \pm 0.79 [0]	69.29 \pm 0.74 [4]
ADMP-GCN (ST) - Oracle	89.43\pm0.19	81.96\pm0.49	97.24\pm0.52	90.13\pm0.36	85.97\pm7.59	79.64\pm0.27

consider the Walk Count centrality for $\ell = 2$, as it effectively captures higher-order structures and provides valuable insights into the organization of complex systems [1].

Using this assumption, we employ node clustering to partition the node set into C clusters, denoted by $\mathcal{P} = \cup_{1 \leq c \leq C} \mathcal{P}_c$. Each cluster $c \in \{1, \dots, C\}$ is assigned a common exit layer $\ell_c \in \{0, \dots, L\}$, determined using nodes excluded from the test set. Validation nodes are utilized for this purpose. (i) Centrality scores are computed for all nodes, considering metrics. (ii) Nodes are ranked based on their centrality scores. (iii) These centrality scores are discretized into C equal-sized buckets to facilitate the clustering process. (iv) The optimal exit layer for each cluster is determined by evaluating the classification accuracy on validation nodes within that cluster, i.e.,

$$\forall c \in \{1, \dots, C\}, \quad \ell_c = \arg \max_{\ell \in \{0, \dots, L\}} \text{Acc} \{p_v^{(\ell)} \in \mathcal{V}_{val} \cap \mathcal{P}_c\}.$$

For the cluster-based layer selection policy, we considered alternative deep learning approaches to predict the optimal exit layer for each node. This high accuracy leads to a significant distribution shift between the predictions for train nodes and those for test nodes. Conversely, validation nodes, which were not utilized during the training of the ADMP-GNN, present a more suitable option for learning the policy due to their unbiased predictions. For the choice of the cluster-based layer selection policy, we could use a variety of deep learning mechanisms to predict the best exit layer for each node. This included generalizing the policy by training neural networks to identify layers that accurately predict node outcomes or by framing the problem as an optimal stopping problem following the framework of [21]. However, these approaches require learning the policy on a set of nodes larger than the test set, which is impractical for node classification tasks where the training and validation node sets available for policy learning are typically smaller than the test set.

4 Experimental Results

4.1 Experimental Setup

Datasets. We use thirteen widely used datasets in the GNN literature. We particularly used the citation networks Cora, CiteSeer, and PubMed [35], the co-authorship networks CS [36], the citation network between Computer Science arXiv papers Ogbn-arxiv [19], the Amazon Computers and Amazon Photo networks [36], the non-homophilous dataset genius [26], and the disassortative datasets Chameleon, Squirrel [32], and Cornell, Texas, Wisconsin from the WebKB dataset [27]. More details and statistics about the used datasets can be found in Appendix A. For the Cora, CiteSeer, and Pubmed datasets, we used the provided train/validation/test

splits. For the remaining datasets, we followed the framework in [27, 32].

Algorithm 1: Sequential Training for ADMP-GNN

Inputs: Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$, number of layers L , node classification loss function \mathcal{L} ,

foreach $t \in \{0, \dots, L\}$ **do**

if $t = 0$ **then**

 1. Set $h_v^{(0)} \leftarrow m_v^{(0)} = x_v$ for all $v \in \mathcal{V}$.

 2. Compute predictions at layer $\ell = 0$, i.e.,

$$\forall v \in \mathcal{V}, \quad p_v^{(0)} = \phi_{\text{Ex}}^{(0)}(h_v^{(0)}).$$

 3. Train the weights of $\phi_{\text{Ex}}^{(0)}$ to minimize $\mathcal{L}(p_v^{(0)})$.

 4. Freeze the gradients of $\phi_{\text{Ex}}^{(0)}$.

else

 1. Use the $\phi_{\text{Ct}}^{(t-1)}$ to update node representations

$$\forall v \in \mathcal{V}, \quad \tilde{h}_v^{(t-1)} = \phi_{\text{Ct}}^{(t-1)}(h_v^{(t-1)}).$$

 2. Aggregate the information for neighbor nodes

$$\forall v \in \mathcal{V}, \quad m_v^{(t)} = \psi^{(t)}(\{\tilde{h}_u^{(t-1)} : u \in \mathcal{N}(v)\}).$$

 3. Compute predictions at layer $\ell = t$, i.e.,

$$\forall v \in \mathcal{V}, \quad p_v^{(t)} = \phi_{\text{Ct}}^{(t-1)}(\tilde{h}_v^{(t-1)}, m_v^{(t)}).$$

 4. Train $\phi_{\text{Ct}}^{(t-1)}$, $\psi^{(t)}$, and $\phi_{\text{Ex}}^{(t)}$ to minimize $\mathcal{L}(p_v^{(t)})$.

 5. Freeze the gradients of $\phi_{\text{Ct}}^{(t-1)}$, $\psi^{(t)}$, and $\phi_{\text{Ex}}^{(t)}$.

end foreach

Baselines. We compare our approach with architectures that combine all the hidden representations of nodes to form a final node representation used for prediction. For each baseline model, we vary the number of layers from 0 to 5, and we report in Table 3 the performance of the best number of layers with respect to the test set. (i) This includes *Jumping knowledge*, which combines the nodes representation of all layers using an aggregation layer, e.g., MaxPooling (JKMaxPool), Concatenation (JK-Concat), or LSTM-attention (JK-LSTM) [42]. (ii) Residuals-GCNI uses an initial residual connection and an identity mapping for each layer. The initial residual connection ensures that the final representation of each node retains at least a fraction of α from the input layer [7]. (iii) GPR-GCN combines adaptive generalized PageRank (GPR) scheme with GNNs [8]. (iv) Ada-GCN, which proposes an RNN-like deep GNN architecture by incorporating AdaBoost to combine the layers [37]. In contrast to the baselines, which aggregate the layers and determine the best performance by performing a grid search over the number of layers L , we fix the maximum number of layers

Table 3: Classification accuracy (\pm standard deviation) on different node classification datasets for the baselines based on the GCN backbone. The higher the accuracy (in %), the better the model. Highlighted are the first, second best results. OOM means *Out of memory*.

Model	Cora	CiteSeer	CS	PubMed	Genuis	Ogbn-arxiv
JKNET-CAT	79.52 \pm 1.16 [2]	69.69 \pm 0.05 [1]	91.23 \pm 1.26 [1]	77.63 \pm 0.59 [2]	81.46\pm0.10 [2]	68.54 \pm 0.57 [5]
JKNET-MAX	75.67 \pm 0.18 [1]	70.12 \pm 0.04 [1]	90.55 \pm 0.74 [1]	75.10 \pm 2.41 [2]	80.13 \pm 0.29 [1]	56.66 \pm 0.31 [5]
JKNET-LSTM	78.95 \pm 0.62 [0]	65.83 \pm 1.27 [0]	90.17 \pm 1.41 [2]	77.73 \pm 0.67 [0]	OOM	OOM
Residuals - GCNII	76.84 \pm 0.20 [1]	69.72 \pm 0.06 [1]	90.84 \pm 1.35 [1]	77.82 \pm 0.44 [2]	<u>81.36\pm1.13 [2]</u>	61.48 \pm 3.10 [4]
AdaGCN	75.08 \pm 0.27 [1]	69.58 \pm 0.19 [1]	89.62 \pm 0.51 [0]	76.40 \pm 0.12 [5]	79.85 \pm 0.00 [0]	22.06 \pm 1.67 [5]
GPR-GNN	79.91 \pm 0.43 [2]	69.21 \pm 0.81 [2]	91.42 \pm 1.12 [2]	79.00 \pm 0.39 [2]	81.04 \pm 0.41 [2]	68.03 \pm 0.23 [3]
GCN	80.78 \pm 0.72 [2]	71.25 \pm 0.72 [2]	92.20\pm0.00 [1]	79.32\pm0.41 [2]	80.76 \pm 1.05 [0]	64.37 \pm 0.43 [2]
ADMP-GCN	<u>81.22\pm0.36 [5]</u>	71.33\pm0.40 [2]	<u>91.49\pm0.66 [2]</u>	<u>79.02\pm0.21 [2]</u>	80.17 \pm 0.79 [0]	69.29 \pm 0.74 [4]
ADMP-GCN w/ <i>Degree</i>	81.03 \pm 0.53	71.10 \pm 0.50	91.26 \pm 0.59	78.71 \pm 0.39	80.73 \pm 1.00	<u>69.59\pm0.28</u>
ADMP-GCN w/ <i>k-core</i>	81.19\pm0.40	<u>71.27\pm0.53</u>	91.29 \pm 0.68	78.73 \pm 0.41	80.73 \pm 1.00	69.55 \pm 0.33
ADMP-GCN w/ <i>Walk Count</i>	81.14 \pm 0.41	71.14 \pm 0.50	91.19 \pm 0.64	78.64 \pm 0.60	80.68 \pm 0.97	69.55 \pm 0.34
ADMP-GCN w/ <i>PageRank</i>	81.05 \pm 0.46	71.00 \pm 0.29	91.09 \pm 0.99	78.69 \pm 0.56	81.12 \pm 1.41	69.60\pm0.29

Table 4: Classification accuracy (\pm standard deviation) on different node classification datasets for the baselines based on the GIN backbone. The higher the accuracy (in %), the better the model. Highlighted are the first, second best results. OOM means *Out of memory*.

Model	Cora	CiteSeer	CS	PubMed	Genuis	Ogbn-arxiv
JKNET-CAT	77.94 \pm 0.67 [2]	64.82 \pm 0.04 [1]	89.26 \pm 1.18 [1]	75.89 \pm 2.50 [2]	OOM	60.23 \pm 0.37 [1]
JKNET-MAX	77.47 \pm 0.81 [1]	64.96 \pm 2.08 [2]	87.40 \pm 2.11 [0]	76.21 \pm 1.73 [0]	OOM	51.84 \pm 4.01 [0]
JKNET-LSTM	77.39 \pm 1.39 [2]	64.52 \pm 2.02 [1]	87.27 \pm 3.66 [3]	75.90 \pm 1.63 [0]	OOM	OOM
GPR-GIN	76.83 \pm 1.22 [2]	<u>66.43\pm1.15 [2]</u>	88.15 \pm 1.53 [5]	77.27 \pm 0.87 [2]	80.82 \pm 0.42 [3]	63.05\pm0.44 [2]
GIN	77.73 \pm 0.99 [2]	65.23 \pm 1.45 [2]	90.29 \pm 0.99 [1]	76.05 \pm 1.14 [2]	80.78 \pm 1.03 [0]	60.70 \pm 0.15 [1]
ADMP-GIN	78.07 \pm 0.68 [2]	65.41 \pm 1.91 [2]	90.82 \pm 1.15 [1]	76.46 \pm 1.04 [4]	80.47 \pm 0.91 [0]	60.85 \pm 0.01 [1]
ADMP-GIN w/ <i>Degree</i>	<u>78.12\pm0.70</u>	66.82\pm0.89	90.70 \pm 0.79	76.07 \pm 1.27	<u>81.68\pm0.70</u>	60.85 \pm 0.01
ADMP-GIN w/ <i>k-core</i>	78.10 \pm 0.64	66.36 \pm 1.03	90.85\pm0.93	75.93 \pm 1.13	81.48 \pm 0.64	<u>60.85\pm0.01</u>
ADMP-GIN w/ <i>Walk Count</i>	78.19\pm0.68	65.79 \pm 0.81	<u>90.77\pm0.91</u>	<u>76.72\pm0.76</u>	81.21 \pm 0.58	60.85 \pm 0.01
ADMP-GIN w/ <i>PageRank</i>	77.78 \pm 0.83	67.08 \pm 0.51	90.72 \pm 0.91	76.63 \pm 0.87	81.89\pm1.04	60.85 \pm 0.01

to $L = 5$ for ADMP-GNN. This comparison gives a strong advantage to the baselines, as their results are optimized for each dataset, while ADMP-GNN’s performance is evaluated under a fixed and consistent setup.

Implementation Details. We train all the models using the Adam optimizer [23] and the same hyperparameters. The GNN hyperparameters in each dataset were performed using a Grid search on the classical GCN; we detail the values of these hyperparameters in Table 7 of Appendix D. To account for the impact of random initialization, each experiment was repeated 10 times, and the mean and standard deviation of the results were reported. The experiments have been run on both an NVIDIA A100 GPU.

4.2 Experimental Results

Through extensive experiments on multiple datasets, we can better understand the scenarios in which ADMP-GNN proves to be

effective. As observed in Tables 3, 4, a comparison between ADMP-GCN and ADMP-GIN against their respective baselines, GCN and GIN, demonstrates consistently higher accuracy for most datasets. Regarding the centrality-based layer selection policy, it becomes clear that this policy is particularly efficient when graphs exhibit a wide range of local density and centrality among nodes. Most importantly, no impactful drop in accuracy was observed with ADMP-GNN. It is important to note that the suitability of a given centrality metric can vary significantly across datasets. For instance, while PageRank is widely used, its distribution often forms a peak near zero due to the normalization constraint where the sum of PageRank scores across nodes equals 1, making it challenging to form meaningful clusters. In contrast, other centrality metrics like k -core and Walk Count are more adaptable in datasets where clear clustering patterns emerge. For example, in datasets like Cora, Ogbn, and Photo, the k -core centrality effectively highlights clusters, enabling an adaptive depth policy. For datasets like Texas and Wisconsin, Walk Count effectively captures the structural diversity necessary

for cluster-based layer selection policy. These observations emphasize the importance of selecting appropriate centrality metrics tailored to the specific properties of a dataset to maximize the effectiveness of ADMP-GNN.

5 Conclusion

In this work, we propose ADMP-GNN, a novel adaption of message passing neural networks that enables to make predictions for each node at every layer. Additionally, we propose a sequential training approach designed to achieve performance comparable to training multiple GNNs independently in a single-task setting. Our empirical analysis highlights the importance of node-specific depth in GNNs to effectively capture the unique characteristics and computational requirements of each node. Determining the optimal number of message-passing layers for each node is a challenging task, influenced by factors such as the complexity and connectivity of graph structures and the variability introduced by node features. Through our experiments, we identified that node centrality can serve as a useful indicator for determining the optimal layer for each node. To this end, we heuristically learn a layer selection policy using a set of validation nodes, which is then generalized to test nodes. Extensive experiments across multiple datasets, particularly those characterized by a diversity in local structural properties, demonstrate that ADMP-GNN significantly enhances the prediction accuracy of GNNs, offering an effective solution to address the challenges of layer selection and node-specific learning.

A Dataset Statistics

Characteristics and information about the datasets utilized in the node classification part of the study are presented in Table 5, which provides statistics about each dataset, and highlights the variety in their features, size, and edge homophily,

Table 5: Statistics of the node classification datasets used in our experiments.

Dataset	#Features	#Nodes	#Edges	#Classes	Edge Homophily
Cora	1,433	2,708	5,208	7	0.809
CiteSeer	3,703	3,327	4,552	6	0.735
PubMed	500	19,717	44,338	3	0.802
CS	6,805	18,333	81,894	15	0.808
Genuis	12	421,961	984,979	5	0.618
Ogbn-arxiv	128	169,343	2,315,598	40	0.654

B Node-Specific Depth Analysis in Graph Neural Networks

We generate synthetic graphs of size $N = 5,000$. We select nodes belonging to sparse or dense regions in the original graph based on their core number. We consider only nodes with labels that are sufficiently present in both dense and sparse region. Last, we randomly select nodes, all by keeping the label distribution similar in both sparse and dense subgraphs. Yellow points indicate edges, while purple points represent non-edges. Notably, there are variations in edge density across different blocks: the first block is characterized

by extreme sparsity, whereas the second block exhibits a much denser structure.

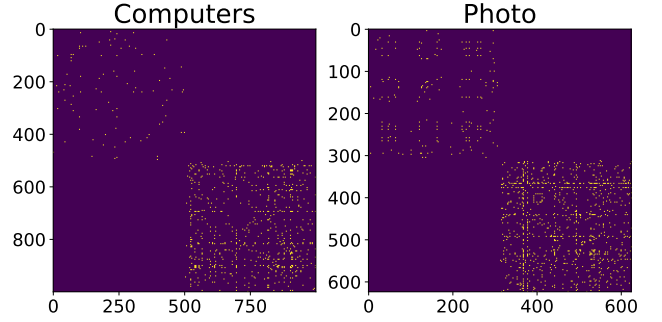


Figure 3: The adjacency matrix of the synthetic graphs extracted from the real graphs Computers and Photo.

C Time Complexity

Understanding the time complexity of the ADMP-GCN model is crucial for evaluating its practical efficiency and scalability. Table 6 reports the average training time, measured in seconds, for two distinct settings: ALM and ST.

Table 6: The average time needed for each training setting for different datasets.

Model	Cora	squirrel	chameleon	Computers	Photo	Ogbn-arxiv
ADMP-GCN ALM	14	36	15	51	26	266
ADMP-GCN ST	32	88	40	175	87	882

D Hyperparameter Configurations

For a more balanced comparison, however, we use the same training procedure for all the models. The hyperparameters in each dataset were performed using a Grid search on the classical GCN over the following search space:

- Hidden size: [16, 32, 64, 128, 256, 512],
- Learning rate: [0.1, 0.01, 0.001],
- Dropout probability: [16, 32, 64, 128, 256, 512],

The number of layers was fixed to 2. The optimal hyperparameters can be found in Table 7.

Table 7: Hyperparameters used in our experiments.

Dataset	Hidden Size	Learning Rate	Dropout Probability
Cora	64	0.01	0.8
CiteSeer	64	0.01	0.4
PubMed	64	0.01	0.2
CS	512	0.01	0.4
Genuis	512	0.01	0.2
Ogbn-Arxiv	512	0.01	0.5

E Supplementary Results of ADMP-GIN

We replicate the experiments from the main papers described in Section 3, using the GIN backbone. The results for ADMP-GIN are presented in Tables 8 and 9.

F GenAI Usage Disclosure

No generative AI tools were used in the preparation of this research, including in the code, data, analysis, or writing of this paper.

References

- [1] Austin R Benson, David F Gleich, and Jure Leskovec. 2016. Higher-order organization of complex networks. *Science* 353, 6295 (2016), 163–166.
- [2] Tolga Bolukbasi, Joseph Wang, Ofer Dekel, and Venkatesh Saligrama. 2017. Adaptive neural networks for efficient inference. In *International Conference on Machine Learning*. PMLR, 527–536.
- [3] Stefan Bornholdt and Heinz Georg Schuster. 2001. Handbook of graphs and networks. *From Genome to the Internet*, Willey-VCH (2003 Weinheim) (2001).
- [4] Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems* 30, 1-7 (1998), 107–117.
- [5] Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Congrui Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, et al. 2020. Spectral temporal graph neural network for multivariate time-series forecasting. *Advances in neural information processing systems* 33 (2020), 17766–17778.
- [6] Jhon A. Castro-Correa, Jhony H. Giraldo, Mohsen Badiéy, and Fragkiskos D. Malliaros. 2024. Gegenbauer Graph Neural Networks for Time-Varying Signal Reconstruction. *IEEE Transactions on Neural Networks and Learning Systems* 35, 9 (2024), 11734–11745.
- [7] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and deep graph convolutional networks. In *International conference on machine learning*. PMLR, 1725–1735.
- [8] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. 2020. Adaptive universal generalized pagerank graph neural network. *arXiv preprint arXiv:2006.07988* (2020).
- [9] Gabriele Corso, Bowen Jing, Regina Barzilay, Tommi Jaakkola, et al. 2023. DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking. In *International Conference on Learning Representations (ICLR 2023)*.
- [10] Alexandre Duval, Victor Schmidt, Alex Hernández-García, Santiago Miret, Fragkiskos D. Malliaros, Yoshua Bengio, and David Rolnick. 2023. FAENet: Frame Averaging Equivariant GNN for Materials Modeling. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR, 9013–9033.
- [11] Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. 2020. Depth-Adaptive Transformer. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=SJg7KhVKPH>
- [12] Moshe Eliasof, Beatrice Bevilacqua, Carola-Bibiane Schönlieb, and Haggai Maron. 2024. GRANOLA: Adaptive Normalization for Graph Neural Networks. *arXiv preprint arXiv:2404.13344* (2024).
- [13] Federico Errica, Henrik Christiansen, Viktor Zaverkin, Takashi Maruyama, Mathias Niepert, and Francesco Alesiani. 2023. Adaptive Message Passing: A General Framework to Mitigate Oversmoothing, Oversquashing, and Underreaching. *arXiv preprint arXiv:2312.16560* (2023).
- [14] Lukas Faber and Roger Wattenhofer. 2024. GwAC: GNNs with Asynchronous Communication. In *Learning on Graphs Conference*. PMLR, 8–1.
- [15] Ben Finkelstein, Xingyue Huang, Michael M. Bronstein, and Ismail Ilkan Ceylan. 2024. Cooperative Graph Neural Networks. <https://openreview.net/forum?id=T0FuEDnODP>
- [16] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *International conference on machine learning*. PMLR, 1263–1272.
- [17] Jhony H. Giraldo, Konstantinos Skianis, Thierry Bouwmans, and Fragkiskos D. Malliaros. 2023. On the Trade-off between Over-smoothing and Over-squashing in Deep Graph Neural Networks. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*. 566–576.
- [18] Alex Graves. 2016. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983* (2016).
- [19] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems* 33 (2020), 22118–22133.
- [20] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. 2016. Deep networks with stochastic depth. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*. Springer, 646–661.
- [21] Côme Huré, Huyén Pham, Achref Bachouch, and Nicolas Langrené. 2021. Deep neural networks algorithms for stochastic control problems on finite horizon: convergence analysis. *SIAM J. Numer. Anal.* 59, 1 (2021), 525–557.
- [22] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. 2016. Dynamic filter networks. *Advances in neural information processing systems* 29 (2016).
- [23] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. doi:10.48550/ARXIV.1412.6980
- [24] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
- [25] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. 2017. Pruning Filters for Efficient ConvNets. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJqFGTslg>
- [26] Derek Lim and Austin R Benson. 2021. Expertise and dynamics within crowd-sourced musical knowledge curation: A case study of the genius platform. In

Table 8: Comparison of ADMP-GIN training paradigms ALM and ST. These paradigms are also compared to the single-task training setting to evaluate which approach most closely mimics the classical GIN under single-task training. The best results for each dataset are bolded.

# Layers	Training Paradigm	Model	Cora	CiteSeer	CS	PubMed	Genius	Ogbn-arxiv
①	Single-task	GIN	56.40±0.06	57.14±0.09	87.17±1.41	72.49±0.08	80.78±1.03	48.87±0.04
	Multi-task	ADMP-GIN (ALM)	58.00±0.00	61.50±0.00	86.36±0.78	73.20±0.00	79.95±0.10	36.49±0.19
		ADMP-GIN (ST)	56.38±0.04	57.17±0.08	87.41±0.95	72.47±0.14	80.47±0.91	48.87±0.04
①	Single-task	GIN	75.17±0.09	64.79±0.03	90.29±0.99	74.97±0.11	78.42±4.95	60.90±0.15
	Multi-task	ADMP-GIN (ALM)	74.50±0.00	66.50±0.00	88.66±1.18	75.40±0.00	72.07±17.44	59.43±0.73
		ADMP-GIN (ST)	75.07±0.06	64.80±0.00	90.82±1.15	75.00±0.00	77.01±12.06	60.85±0.01
②	Single-task	GIN	77.73±0.99	65.23±1.45	87.93±0.71	76.05±1.14	78.89±0.48	16.23±9.60
	Multi-task	ADMP-GIN (ALM)	63.53±4.80	60.68±2.09	12.90±8.05	72.20±4.11	79.60±0.51	47.63±3.20
		ADMP-GIN (ST)	78.07±0.68	65.41±1.91	87.47±2.18	76.04±0.88	72.99±13.48	10.13±8.00
③	Single-task	GIN	74.40±1.14	60.81±2.20	82.13±2.24	74.97±1.69	52.22±28.47	6.00±0.17
	Multi-task	ADMP-GIN (ALM)	66.99±2.85	59.57±2.57	17.69±12.02	74.15±2.00	68.00±23.98	27.51±16.25
		ADMP-GIN (ST)	76.28±1.11	65.13±1.00	83.60±3.72	76.21±1.75	80.04±0.09	13.74±9.66
④	Single-task	GIN	67.68±4.07	57.54±2.92	47.37±17.10	73.62±1.63	80.05±0.10	6.07±0.21
	Multi-task	ADMP-GIN (ALM)	68.78±4.30	60.12±2.36	21.07±14.15	74.20±2.57	79.36±1.18	15.23±11.85
		ADMP-GIN (ST)	74.94±1.58	65.21±1.54	81.33±2.15	76.46±1.04	80.04±0.09	16.02±10.78
⑤	Single-task	GIN	32.48±10.47	54.76±2.32	18.56±11.32	67.98±5.80	80.05±0.10	6.19±0.39
	Multi-task	ADMP-GIN (ALM)	67.46±4.34	59.74±1.95	29.16±13.84	73.94±2.42	79.99±0.10	14.61±12.11
		ADMP-GIN (ST)	71.34±2.09	63.89±1.39	79.10±1.84	75.75±1.09	79.57±1.40	15.72±12.13

Table 9: Comparison of highest accuracy (\pm standard deviation) for GIN and ADMP-GIN ST, with the layer achieving the best accuracy indicated in brackets. The final row shows the *Oracle Accuracy* for ADMP-GIN ST. Best results for each dataset are bolded.

Model	Cora	CiteSeer	CS	PubMed	Genius	Ogbn-arxiv
GIN	77.73±0.99 [2]	65.23±1.45 [2]	90.29±0.99 [1]	76.05±1.14 [2]	80.78±1.03 [0]	60.90±0.15 [1]
ADMP-GIN (ST)	78.07±0.68 [2]	65.41±1.91 [2]	90.82±1.15 [1]	76.46±1.04 [4]	80.47±0.91 [0]	60.85±0.01 [1]
ADMP-GIN ST - Oracle	90.76±0.27	81.87±0.63	97.64±0.18	92.73±0.54	92.07±6.13	71.23±2.69

Proceedings of the International AAAI Conference on Web and Social Media, Vol. 15. 373–384.

- [27] Derek Lim, Felix Matthew Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Prasad Bhalerao, and Ser-Nam Lim. 2021. Large Scale Learning on Non-Homophilous Graphs: New Benchmarks and Strong Simple Methods. In *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (Eds.). <https://openreview.net/forum?id=DfGu8WwT0d>
- [28] Juncheng Liu, Kenji Kawaguchi, Bryan Hooi, Yiwei Wang, and Xiaokui Xiao. 2021. Eignn: Efficient infinite-depth graph neural networks. *Advances in Neural Information Processing Systems* 34 (2021), 18762–18773.
- [29] Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. 2022. Revisiting heterophily for graph neural networks. *Advances in neural information processing systems* 35 (2022), 1362–1375.
- [30] Fragkiskos D. Malliaros, Christos Giatsidis, Apostolos N. Papadopoulos, and Michalis Vazirgiannis. 2020. The core decomposition of networks: theory, algorithms and applications. *Vldb j.* 29, 1 (2020), 61–92.
- [31] Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. 2022. Recipe for a general, powerful, scalable graph transformer. *Advances in Neural Information Processing Systems* 35 (2022), 14501–14515.
- [32] Benedek Rozemberczki, Carl Allen, and Rik Sarkar. 2021. Multi-scale attributed node embedding. *Journal of Complex Networks* 9, 2 (2021), cnab014.
- [33] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. 2017. Dynamic routing between capsules. *Advances in neural information processing systems* 30 (2017).
- [34] Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q. Tran, Yi Tay, and Donald Metzler. 2022. Confident Adaptive Language Modeling. In *Advances in Neural Information Processing Systems*, Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (Eds.). <https://openreview.net/forum?id=uLYc4L3C81A>
- [35] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. 2008. Collective Classification in Network Data. *AI Magazine* 29, 3 (Sep. 2008), 93. doi:10.1609/aimag.v29i3.2157
- [36] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Pitfalls of Graph Neural Network Evaluation. doi:10.48550/ARXIV.1811.05868
- [37] Ke Sun, Zhanxing Zhu, and Zhouchen Lin. 2021. Ada[GCN]: Adaboosting Graph Convolutional Networks into Deep Models. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=QkRbdiiEjM>
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *CoRR* abs/1706.03762 (2017). arXiv:1706.03762 <http://arxiv.org/abs/1706.03762>

- [39] Yulin Wang, Kangchen Lv, Rui Huang, Shiji Song, Le Yang, and Gao Huang. 2020. Glance and Focus: a Dynamic Approach to Reducing Spatial Redundancy in Image Classification. *CoRR* abs/2010.05300 (2020). arXiv:2010.05300 <https://arxiv.org/abs/2010.05300>
- [40] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In *7th International Conference on Learning Representations*.
- [41] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation Learning on Graphs with Jumping Knowledge Networks. *CoRR* abs/1806.03536 (2018). arXiv:1806.03536 <http://arxiv.org/abs/1806.03536>
- [42] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation Learning on Graphs with Jumping Knowledge Networks. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 5453–5462. <https://proceedings.mlr.press/v80/xu18c.html>
- [43] Hanqing Zeng, Muhan Zhang, Yinglong Xia, Ajitesh Srivastava, Rajgopal Kannan, Viktor Prasanna, Long Jin, Andrey Malevich, and Ren Chen. 2020. Deep graph neural networks with shallow subgraph samplers. (2020).
- [44] Wentao Zhang, Zeang Sheng, Yuezihan Jiang, Yikuan Xia, Jun Gao, Zhi Yang, and Bin Cui. 2021. Evaluating deep graph neural networks. *arXiv preprint arXiv:2108.00955* (2021).
- [45] Lingxiao Zhao and Leman Akoglu. 2020. PairNorm: Tackling Oversmoothing in GNNs. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rkecl1rtwB>
- [46] Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian J. McAuley, Ke Xu, and Furu Wei. 2020. BERT Loses Patience: Fast and Robust Inference with Early Exit. *CoRR* abs/2006.04152 (2020). arXiv:2006.04152 <https://arxiv.org/abs/2006.04152>