

# HIERARCHICAL GRAPH NEURAL NETS CAN CAPTURE LONG-RANGE INTERACTIONS

Ladislav Rampášek    Guy Wolf

Université de Montréal, Dept. of Math. & Stat.; Mila - Quebec AI Institute, Montreal, QC, Canada

## ABSTRACT

Graph neural networks (GNNs) based on message passing between neighboring nodes are known to be insufficient for capturing long-range interactions in graphs. In this paper we study hierarchical message passing models that leverage a multi-resolution representation of a given graph. This facilitates learning of features that span large receptive fields without loss of local information, an aspect not studied in preceding work on hierarchical GNNs. We introduce Hierarchical Graph Net (HGNet), which for any two connected nodes guarantees existence of message-passing paths of at most logarithmic length w.r.t. the input graph size. Yet, under mild assumptions, its internal hierarchy maintains asymptotic size equivalent to that of the input graph. We observe that our HGNet outperforms conventional stacking of GCN layers particularly in molecular property prediction benchmarks. Finally, we propose two benchmarking tasks designed to elucidate capability of GNNs to leverage long-range interactions in graphs.

**Index Terms**— graph neural networks, hierarchical message passing, long range interactions

## 1. INTRODUCTION

Graph neural networks (GNNs), and the field of geometric deep learning, have seen rapid development in recent years [1, 2] and have attained popularity in various fields involving graph and network structures. Prominent examples of GNN applications include molecular property prediction, physical systems simulation, combinatorial optimization, or interaction detection in images and text. Many of the current GNN designs are based on the principle of neural message-passing [3], where information is iteratively passed between neighboring nodes along existing edges. However, this paradigm is known to suffer from several deficiencies, including theoretical limits of their representational capacity [4] and observed limitations of their information propagation over graphs [5, 6, 7].

Two of the most prominent deficiencies of GNNs are known as *oversquashing* and *oversmoothing*. Information

*oversquashing* refers to the exponential growth in the amount of information that has to be encoded by the network with each message-passing iteration, which rapidly grows beyond the capacity of a fixed hidden-layer representation [5]. Signal *oversmoothing* refers to the tendency of node representations to converge to local averages [6], which can also be observed in graph convolutional networks implementing low pass filtering over the graph [7]. A significant repercussion of these phenomena is that they limit the ability of most GNN architectures to represent long-range interactions (LRIs) in graphs. Namely, they struggle in capturing dependencies between distant nodes, even when these have potentially significant impact on output prediction or appropriate internal feature extraction towards it. Capturing LRIs typically requires the number of GNN layers (i.e., implementing individual message passing steps) to be proportional to the diameter of the graph, which in turn exacerbates the oversquashing of massive amount of information and the oversmoothing that tends towards averaging over wide regions of the graph, if not the entire graph.

In this paper, we study the utilization of multiscale hierarchical meta-structures to enhance message passing in GNNs and facilitate capturing of LRIs. By leveraging hierarchical message passing between nodes, our Hierarchical Graph Net (HGNet) architecture can propagate information within  $O(\log |V(G)|)$  steps instead of  $O(\text{diam}(G))$ , leading to particular improvements for sparse graphs with large diameters.

We note that a few works have recently proposed related approaches using hierarchical constructions, namely g-U-Net [8] and GXN [9]. g-U-Net employs a similarity-based top-k pooling called gPool for hierarchical construction over which it implements bottom-up and simple top-down message passing. GXN introduced mutual information based pooling (VIPool) together with a more complex cross-level message passing. Next, MGKN [10] introduced multi-resolution GNN with V-cycle algorithm specifically for learning solutions operators to PDEs. Broadly related are also differentiable pooling methods such as DiffPool [11], EdgePool [12], or GraphZoom [13]. However, these do not employ two-directional hierarchical message passing.

While LRIs are widely accepted as being important for both theoretical studies and in practice, most benchmarks used to empirically validate GNN models do not clearly exhibit this property. Out of these, the importance of LRIs is perhaps best

This work was partially funded by CIFAR AI Chair [G.W.] and IVADO grant PRF-2019-3583139727. The content here is solely the responsibility of the authors and does not necessarily represent the official views of the funding agencies. Code: <https://github.com/rampasek/HGNet>. Correspondence to: {ladislav.rampasek,wolfguy}@mila.quebec

justified in biochemistry datasets, where the 2D structure of proteins and molecules is used as their graph representation. However, edges of such graphs do not encode 3D forces and global properties, leaving it up to the model to learn to recognize such LRIs. Several highly specialized models have been proposed for molecular data, but these are typically not applicable to other domains, which also hinders analysis of their modeling improvements towards particularly capturing LRIs. Therefore, in our experiments we primarily focus on quantifying the benefit of using a hierarchical structure compared to the standard practice of GNN layer stacking. We also introduce two benchmarking tasks designed to elucidate capability of general-purpose GNNs to leverage LRIs. Here, we show hierarchical models outperform their standard GNN counterparts when their hierarchical graph construction matches well with the original graph structure and the prediction task, while uncovering related limitations of gPool in g-U-Net.

## 2. HIERARCHICAL GRAPH NET

To build a hierarchical message passing model, we need to construct a hierarchical graph representation and define an inter- and intra-level message passing mechanism.

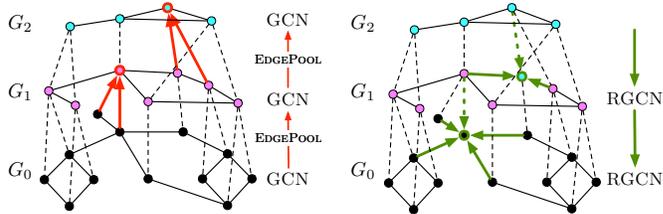
### 2.1. Graph coarsening for hierarchical representation

Building a hierarchical representation principally involves iterative application of graph coarsening and pooling operations. Graph coarsening computes a mapping from nodes of a starting graph  $G_\ell$  onto nodes of a new smaller graph  $G_{\ell+1}$ , while the pooling step computes node and edge features of  $G_{\ell+1}$  from  $G_\ell$ . Here we explore two different approaches: EdgePool [12] and the Louvain method for community detection [14].

**EdgePool** [12] is a method based on the principle of edge contractions. First, the raw score of an edge  $e_{u,v} = (u, v)$  is obtained by a linear combination of respective node features  $x_u$  and  $x_v$ :  $r_{u,v} = W(x_u || x_v) + b$ . Raw scores of edges incident to a node  $u$  are then normalized as  $0.5 + \text{softmax}_{v \in \mathcal{N}(u)} r_{u,v}$  to obtain the final edge scores  $s_{u,v}$ . Finally, a maximal set of edges is greedily selected according to their scores  $s_{a,b}$  and then contracted to create a new graph  $G_{\ell+1}$  from  $G_\ell$ , while nodes in  $G_\ell$  that were not merged are carried forward to  $G_{\ell+1}$ . Two nodes  $a, b$  in  $G_{\ell+1}$  are then connected by an edge iff there exist two nodes in  $G_\ell$  the  $a, b$  were constructed from that had been adjacent in  $G_\ell$ .

Contraction of an edge  $(u^{(\ell)}, v^{(\ell)}) \in G_\ell$  results in a new node  $w^{(\ell+1)}$  with features  $x_w^{(\ell+1)} := s_{u,v}^{(\ell)}(x_u^{(\ell)} + x_v^{(\ell)})$ . Multiplying the new node features by the edge score facilitates gradient-based learning of the scoring function, which would otherwise be independent of the final objective function.

**Louvain method** for community detection [14] is a heuristic method based on greedy maximization of modularity score of each community. It is an  $O(N \log N)$  algorithm without learnable parameters that is deterministic for a fixed random



**Fig. 1.** HGNet with two hierarchical levels over an original graph  $G_0$  of 12 vertices (in black) and 14 edges. The dashed lines represent inter-level edges. (left) Two levels of EdgePool coarsening, highlighted by red arrows, create the hierarchical structure. A GCN layer is applied before each EdgePool coarsening and at the final coarsest level  $G_2$ . (right) Message passing down the hierarchy is implemented by an RGCN layer at  $G_1$  and then  $G_0$  levels, highlighted by green arrows, where inter-level edges are treated as a distinct edge type.

seed. The Louvain algorithm merges clusters (communities) into a single node and iteratively performs modularity clustering on the condensed graph until the score cannot be improved. The size of the condensed graph cannot be directly controlled, but seems to yield satisfying contraction ratios in practice.

To build a hierarchical meta-graph over a starting graph  $G_0$ , we use average node and edge feature pooling according to the modular communities identified in  $G_\ell$  by the Louvain method to construct the following level  $G_{\ell+1}$ .

### 2.2. Hierarchical message passing in HGNet

Both EdgePool and the Louvain method provide a recipe for construction of a hierarchical graph representation. We propose Hierarchical Graph Network (HGNet) based on either one of these messages (see Figure 1), sharing the same hierarchical message passing approach that we describe next. Our message passing both within and between levels is principally similar to that of g-U-Net. Consider a hierarchical meta-graph with  $L$  levels over some  $G_0$ . The forward propagation in HGNet consists of a computational pass going up the hierarchy and of a pass going down the hierarchy, resulting in the final embedding of each node in  $G_0$ . In the upwards pass we first apply a GCN layer [15] to  $G_\ell$ , starting with  $\ell = 0$ , followed by node and edge pooling according to either EdgePool or the Louvain method to instantiate the next hierarchical level  $G_{\ell+1}$ . This process iterates until the final level  $L$ , at which point no more pooling is done and the downwards pass starts. In this downwards pass we utilize RGCN [16] layers at each  $G_\ell$  level  $\ell \in \{L-1, \dots, 0\}$ , where we add special edges that connect merged nodes in  $G_\ell$  with their respective representatives in  $G_{\ell+1}$  by an edge of unique type.

**Complexity.** We now analyze the asymptotic complexity of our hierarchical meta-graph based on the EdgePool variant. Let us assume that in each round of edge contractions the size of the greedy maximum matching is at least a constant

fraction  $m \geq 2$  of the number of remaining nodes, i.e.,  $\frac{N}{m}$ . Note that  $m = 2$  when the selected set of edges is a perfect matching. That means after the first round there will be  $N \frac{m-1}{m}$  nodes in the next  $G_1$  level. Thus, the total number of nodes in the entire hierarchical structure over a  $G_0$  with  $N$  nodes is  $\sum_{\ell=0}^{\infty} N \left(\frac{m-1}{m}\right)^\ell = mN = O(N)$ , while the number of possible levels is  $\log_{\frac{m-1}{m}} N = O(\log N)$ . This construction therefore guarantees that, if  $G_0$  is connected, the shortest path length between any two nodes is upper-bounded by  $O(\log N)$ .

We can also expect the number of edges in our hierarchical graph to remain asymptotically equal to the number of edges in the input graph  $G_0$ . Assume there are  $E = \Omega(N)$  edges in  $G_0$  out of  $O(N^2)$  possible and that they are uniformly distributed. Then after one round of EdgePool, the number of edges in  $G_1$  is expected to be  $O\left(E \left(\frac{m-1}{m}\right)^2\right)$ , because the number of possible edges in  $G_1$  compared to  $G_0$  has decreased from  $O(N^2)$  to  $O\left(\left[N \frac{m-1}{m}\right]^2\right)$ , i.e., we can expect  $\left(\frac{m-1}{m}\right)^2$  contraction factor for the number of edges. Therefore, we can expect  $\sum_{\ell=0}^{\infty} E \left(\frac{m-1}{m}\right)^{2\ell} = \frac{m^2}{2m-1} E = O(E)$  intra-level edges in total. From the construction of the hierarchy it is also clear that the number of inter-level edges (connecting nodes between adjacent hierarchical levels) is  $O(N)$  as the total number of nodes is  $O(N)$ . Therefore, the total number of edges is expected to remain  $O(E)$ .

Given a deep enough hierarchy and large enough node representation capacity, the final node embeddings can incorporate LRIs from the entire graph  $G_0$ , as well as local information. In the case of EdgePool, the asymptotic complexity of our HGNet remains that of GCN, as even despite our hierarchical graph having up to  $O(\log N)$  hierarchical levels, its size remains asymptotically unchanged under reasonable assumptions. For a standard message passing GNN to theoretically achieve this capability, it is necessary to stack  $O(\text{diam}(G))$  layers, which may be prohibitively expensive.

### 3. RESULTS

In order to evaluate the performance of HGNet, we consider a wide variety of graph data, including transductive node classification and inductive graph-level classification. Our benchmarks include two settings of HGNet (namely, with EdgePool and Louvain hierarchical structures) and six competitive baseline models: GCN [15], GCN+VN (GCN extended with a Virtual Node connected to all other nodes), GAT [17], ChebNet [18], GIN [4], and g-U-Net [8]. The experimental setup is identical for all tested methods. Each method is trained for 200 epochs, followed by a selection of the best model based on the validation performance, and finally performance on the test split is reported. In case of GCN, GCN+VN, GAT, ChebNet and GIN, we always used a stack of 2 layers unless explicitly stated otherwise. In the case of g-U-Net, we reproduced published hyperparameters [8] as closely as possible. For each method we default to 32-dimensional hidden

node representation; other hyperparameters specific to certain tasks or datasets are described in the respective sections. We note that our reproduced g-U-Net results differ from the original publication [8], as there only the best validation set results were reported rather than performance on independent test sets. This erroneous practice had occurred on several occasions in the relatively nascent field of graph deep learning [19].

#### 3.1. Node classification in citation networks

For our first benchmark, we consider semi-supervised node classification on the CiteSeer, Cora and PubMed citation networks [20]. Our HGNet variants are configured with one hierarchical level and g-U-Net with four levels as per published hyperparameters. Citation networks are known to exhibit high homophily [21], i.e., nodes tend to have the same class label as most of their first degree neighbors. First-order message passing GNNs are known to perform well in high-homophily settings [21], which is validated by our experiments presented in Table 1, with the exception of GCN+VN and GIN. All three hierarchical methods (i.e., g-U-Net, HGNet-EdgePool, and HGNet-Louvain) attain very similar results, slightly behind the best performing GAT, GCN, and ChebNet.

The low performance of GCN+VN, a model geared towards capturing global information, and middle-of-the-pack performances of the hierarchical methods can be explained by the high homophily present in the data, and support prior findings [22] showcasing that global graph information is not vital in these datasets. Hence, given similar model capacity and experimental settings, methods favoring local information, such as GAT and GCN, outperform the more sophisticated ones. We conclude that CiteSeer, Cora and PubMed are not directly suitable to test the ability of GNN models to capture global information or LRIs, despite their extensive use and popularity in such benchmarks [8, 9].

#### Resampled citation networks

In an effort to make the prediction tasks of CiteSeer, Cora and PubMed citation networks more suitable for testing the models’ ability to utilize information from farther nodes, we experimented with a specific resampling of their training, validation and test splits. The standard semi-supervised splits [20] follow the same key for each dataset: 20 examples from each class are randomly selected for training, while 500 and 1000 examples are drawn uniformly randomly for the validation and test splits. We used principally the same key, but a different random sampling strategy. Once a node is drawn, we enforced that none of its  $k$ -th degree neighbors is selected for any split. This approach guarantees that a  $k$ -hop neighborhood of each labeled node is “sanitized” of labels. As such, we prevent potential correct-class label imprinting in the representation of these  $k$ -th degree neighbors during the semi-supervised transductive training. For a model to leverage such imprinting benefit of homophily, it has to be able to reach beyond this

**Table 1. Legacy graph benchmarks.** CiteSeer, Cora and PubMed provide only one standard data split, and therefore we show test accuracy averaged over three runs with different random seeds for these datasets. For graph classification tasks (right side of the table) we used 10-fold stratified cross-validation. Shown heatmaps are normalized per dataset (column).

	CiteSeer	Cora	PubMed	COLLAB	IMDB-B	IMDB-M	D&D	NCI1	ENZYMES	PROTEINS
GAT	70.10%	81.97%	76.27%	78.90%	71.60%	50.20%	69.27%	63.70%	34.17%	71.43%
GCN	70.00%	80.90%	77.70%	78.66%	71.30%	49.47%	71.56%	63.82%	29.17%	71.52%
GCN+VN	23.77%	31.67%	42.63%	79.80%	69.50%	47.17%	79.66%	71.47%	40.83%	75.18%
ChebNet	69.40%	79.43%	78.03%	79.48%	73.50%	49.80%	70.29%	67.66%	34.33%	71.61%
GIN	56.33%	71.80%	73.63%	78.90%	71.50%	49.83%	70.97%	72.75%	34.17%	70.69%
g-U-Net	66.60%	80.03%	77.10%	81.25%	73.25%	49.83%	71.82%	64.11%	31.67%	72.72%
<b>HGNet-EdgePool</b>	68.37%	80.60%	76.73%	82.15%	70.75%	49.00%	75.64%	77.13%	43.75%	72.48%
<b>HGNet-Louvain</b>	68.60%	81.03%	77.50%	81.13%	72.20%	50.73%	74.36%	75.06%	39.67%	73.77%

**Table 2. Citation networks with  $k$ -hop sanitized dataset splits.** The reported metric is the average test accuracy over three training runs with different random seeds, while keeping the same resampled splits. Heatmaps are normalized per block given by a dataset and neighborhood size  $k$  combination.

model	k=1		k=2			
	layers=1	layers=2	layers=1	layers=2	layers=3	
CiteSeer	GAT	60.63%	66.37%	54.26%	59.35%	55.87%
	GCN	58.00%	64.23%	52.11%	58.26%	56.96%
	GCN+VN	58.00%	30.70%	52.11%	26.29%	24.18%
	ChebNet	53.10%	62.65%	49.35%	55.05%	55.69%
	GIN	52.20%	54.03%	48.00%	43.11%	40.41%
	g-U-Net	63.03%	60.67%	56.26%	57.32%	56.06%
	<b>HGNet-EdgePool</b>	64.33%	62.33%	57.75%	56.92%	56.49%
	<b>HGNet-Louvain</b>	64.10%	61.70%	56.42%	58.41%	59.86%
Cora	GAT	72.08%	78.50%	64.66%	72.51%	73.37%
	GCN	67.90%	76.71%	61.30%	70.59%	72.03%
	GCN+VN	67.81%	34.25%	61.30%	42.72%	42.24%
	ChebNet	60.80%	71.53%	53.59%	69.83%	72.27%
	GIN	67.03%	67.08%	59.39%	62.07%	60.82%
	g-U-Net	78.08%	77.21%	73.47%	74.43%	73.66%
	<b>HGNet-EdgePool</b>	77.17%	75.33%	71.74%	73.95%	72.13%
	<b>HGNet-Louvain</b>	77.12%	76.85%	72.03%	74.90%	75.38%
PubMed	GAT	70.33%	75.87%	67.37%	72.80%	73.10%
	GCN	69.97%	75.83%	67.03%	72.07%	71.30%
	GCN+VN	69.75%	48.30%	67.03%	42.13%	39.40%
	ChebNet	69.20%	73.70%	66.80%	69.35%	71.35%
	GIN	61.80%	73.47%	63.43%	66.50%	68.23%
	g-U-Net	75.83%	75.73%	71.67%	72.37%	70.63%
	<b>HGNet-EdgePool</b>	77.00%	75.30%	71.80%	72.43%	72.87%
	<b>HGNet-Louvain</b>	76.30%	74.17%	70.67%	75.07%	74.30%

$k$ -hop neighborhood, assuming that the class homophily spans that far in the underlying data.

We experimented with  $k \in \{1, 2\}$  for all 3 citation networks and kept the same hyperparameters from the prior experiments, but varied the number of stacked layers or hierarchy levels, as applicable, for each GNN method. Results averaged over runs with 3 random seeds are shown in Table 2. For  $k = 1$  we see consistent degradation of performance for single-layer GNNs, while even one level of hierarchy provides significant advantage for the hierarchical models. GAT and GCN recover competitive performance given two layers, which allows the models to reach second-order neighborhood with some nodes

that are labeled during training. Hierarchical models however do not benefit from using two levels, as with even just one level their receptive field is already large enough to reach beyond first-order neighborhood of a node. In case of  $k = 2$  we observe similar behavior, but now hierarchical models typically benefit from employing two or three levels. This is particularly true for PubMed, the largest tested dataset. In this scenario we believe we have reached the limit of these datasets in the sense that we do not expect third-degree or further nodes to be consistently of significant relevance. We can see that for most methods the performance is relatively similar between two or three layers. Our resampling approach is fundamentally limited by the strong local homophily present in these citation networks and beyond  $k = 2$  cannot be used to test capability of the models to leverage LRIs.

### 3.2. Graph-level prediction

We now turn our focus to graph-level classification. We start by benchmarking all methods using a set of commonly used datasets: COLLAB, IMDB-BINARY, IMDB-MULTI, D&D, NCI1, ENZYMES, and PROTEINS [23]. In the second part we present a new set of datasets we designed to challenge the GNN methods in learning to recognize a complex set of features. In this section, we use global mean pooling for each method to obtain the graph-level representation from individual nodes of a graph. Using this representation, a graph is finally classified by a 2-layer MLP classifier with 128-dimensional hidden layer.

Our experimental results in common graph-classification datasets are presented in Table 1 (right side). One of our HGNet variants is the best performing method in 4 out of the 7 datasets. GCN+VN performs well on molecular datasets where global information is important, as does HGNet. However, g-U-Net falls behind in this setting, likely due to the nature of top-k pooling in its gPool, which destroys local information and appears to have difficulty extracting complex global features.

### OGB molecular benchmarks

We tested HGNet on two Open Graph Benchmark (OGB) [24] molecular property prediction datasets: *ogbg-molpcba* and

**Table 3. OGB molecular benchmarks.** HGNet results are obtained and presented as per OGB standards, shown is the mean and standard deviation from 10 runs with different random seeds. HGNet models have 1, 2, or 3 levels and otherwise mirror hyperparameters of the OGB baselines that each have 5 layers. The metrics for baselines are from the OGB online leaderboard.

	ogbg-molpcba (average precision)			ogbg-molhiv (ROC-AUC)		
	Test	Validation	#Params	Test	Validation	#Params
GCN	0.2020 $\pm$ 0.0024	0.2059 $\pm$ 0.0033	565,928	0.7606 $\pm$ 0.0097	0.8204 $\pm$ 0.0141	527,701
GCN+VN	0.2424 $\pm$ 0.0034	0.2495 $\pm$ 0.0042	2,017,028	0.7599 $\pm$ 0.0119	0.8384 $\pm$ 0.0091	1,978,801
GIN	0.2266 $\pm$ 0.0028	0.2305 $\pm$ 0.0027	1,923,433	0.7558 $\pm$ 0.0140	0.8232 $\pm$ 0.0090	1,885,206
GIN+VN	0.2703 $\pm$ 0.0023	0.2798 $\pm$ 0.0025	3,374,533	0.7707 $\pm$ 0.0149	0.8479 $\pm$ 0.0068	3,336,306
HGNet-EdgePool 1L	0.2012 $\pm$ 0.0028	0.2060 $\pm$ 0.0027	549,429	0.7648 $\pm$ 0.0154	0.8213 $\pm$ 0.0158	511,202
HGNet-EdgePool 2L	0.2133 $\pm$ 0.0033	0.2204 $\pm$ 0.0020	823,030	0.7670 $\pm$ 0.0140	0.8340 $\pm$ 0.0108	784,803
HGNet-EdgePool 3L	0.2146 $\pm$ 0.0028	0.2199 $\pm$ 0.0036	1,096,631	0.7684 $\pm$ 0.0100	0.8325 $\pm$ 0.0078	1,058,404
HGNet-Louvain 1L	0.1995 $\pm$ 0.0033	0.2078 $\pm$ 0.0030	548,828	0.7648 $\pm$ 0.0093	0.8203 $\pm$ 0.0101	510,601
HGNet-Louvain 2L	0.2095 $\pm$ 0.0030	0.2146 $\pm$ 0.0043	821,828	0.7728 $\pm$ 0.0147	0.8327 $\pm$ 0.0087	783,601
HGNet-Louvain 3L	0.2113 $\pm$ 0.0030	0.2137 $\pm$ 0.0029	1,094,828	0.7738 $\pm$ 0.0266	0.8393 $\pm$ 0.0096	1,056,601

**Table 4. Color-connectivity datasets.** The average test accuracy in 10-fold stratified CV for various depths of the models.

dataset model	16x16 grid			32x32 grid			Minnesota			Euroroad		
	layers=1	layers=2	layers=3	layers=1	layers=2	layers=3	layers=1	layers=2	layers=3	layers=1	layers=2	layers=3
GAT	80.46%	84.41%	89.91%	64.65%	77.68%	79.45%	59.28%	63.92%	67.17%	51.98%	52.47%	53.29%
GCN	52.30%	78.91%	88.21%	50.91%	65.87%	78.61%	54.75%	59.05%	72.22%	49.98%	53.03%	58.66%
GCN+VN	51.97%	76.37%	85.05%	51.22%	63.01%	79.78%	54.40%	58.98%	61.08%	50.24%	52.83%	53.02%
ChebNet	53.17%	83.99%	91.27%	52.93%	74.40%	81.42%	51.47%	73.45%	79.18%	51.34%	60.12%	69.52%
GIN	83.40%	89.87%	94.89%	73.89%	81.30%	85.33%	66.20%	77.87%	83.33%	57.12%	73.35%	90.29%
g-U-Net	85.49%	84.92%	81.41%	54.97%	58.47%	56.50%	57.32%	57.03%	57.45%	50.00%	50.00%	50.00%
HGNet-EdgePool	88.18%	92.50%	93.53%	80.31%	84.18%	86.43%	71.55%	81.80%	82.43%	70.49%	88.42%	95.95%
HGNet-Louvain	86.05%	89.88%	91.91%	77.19%	79.45%	84.76%	65.62%	70.63%	74.25%	56.82%	73.63%	90.97%

*ogbg-molhiv*. For our HGNet we used the same experimental setup and GCN layer implementation as provided by OGB. Both EdgePool and Louvain versions of HGNet with 2 hierarchical levels (2L), composed of 3 GCN and 2 RGCN-like layers, outperform GCN with 5 layers (see Table 3). Employing a hierarchical meta-graph is more powerful than stacking the same number of layers. We note that adding global readouts via Virtual Node is remarkably beneficial in *ogbg-molpcba*, albeit at the cost of many additional parameters.

### Color-connectivity task

Open Graph Benchmark and other recent initiatives are increasing the bar for GNN benchmarking, as many established benchmarking datasets are too small or too simple to adequately test the expressive power of new GNN methods. However, the motivation to include a new dataset in a suite is typically based on the interest in a particular application domain and the scale of the dataset. Unfortunately, none of the existing benchmarks provably require the capture of LRIs for significant performance gain. This issue was not realized in the benchmarking of prior hierarchical methods [8, 9], except [25] that proposed shortest path prediction task in random graphs. Here we propose to employ a task not used for GNN benchmarking before –

classifying the connectivity of same colored nodes in graphs of varying topology. Our color-connectivity datasets are created by taking a graph and randomly coloring half of its nodes one color, e.g., red, and the other nodes blue, such that the red nodes either create a single connected island or two disjoint islands. The binary classification task is then distinguishing between these two cases. The node colorings were sampled by running two red-coloring random walks starting from two random nodes. We used 16x16 and 32x32 2D grids, as well as the Euroroad and Minnesota road networks [26] for the underlying graph topology. For each, we sampled a balanced set of 15,000 examples, except for Minnesota network for which we generated 6,000 examples due to memory constraints. Solving this task requires combination of local and long-range information, while a global readout, e.g., via Virtual Node, is expected to be unsatisfactory.

HGNet-EdgePool is the single best method in this suite of benchmarks (Table 4). Given the nature of the data, we observe a large difference in how suitable are the hierarchical graphs created by different approaches. In particular, gPool of g-U-Net fails to facilitate the learning process on large graphs. Next, global readout via Virtual Nodes in the GCN+VN model does not provide any improvement over the standard GCN, as evidently it is not able to capture complex features. On the

other hand, we see that the ChebNet and GIN models perform well. ChebNet can learn filters that have large receptive field in graph space, which is important in this case. We suspect that GIN is powerful enough to learn local heuristics GCN and GAT fail to, which warrants further investigation.

#### 4. CONCLUSION

Across many datasets, we saw hierarchical models outperform their standard GNN counterparts when construction of the hierarchical graph (its inductive bias) matches well with the graph structure and prediction task. We have not compared to methods highly specialized for a particular tasks, e.g., molecular property prediction, but rather focused on elucidating the effect of using a hierarchical structure compared to the standard approach of stacking GNN layers. Further research remains to be done in terms of exploring combinations of various pooling approaches, hierarchical message passing algorithms and utilization of, e.g., GIN layers instead of GCN. Our proposed color-connectivity task requires complex graph processing to which most existing message-passing GNNs do not scale. These datasets can serve as a common-sense validation for new and more powerful methods. Our testbed datasets can still be improved, as the node features are minimal and recognition of particular topological patterns (e.g., rings or other subgraphs) is not needed to solve the current task. Nevertheless, it represents a significant step forward in terms of understanding and benchmarking more complex graph neural networks.

**Acknowledgments:** The authors would like to thank William L. Hamilton for insightful discussions and Semih Cantürk for help with proofreading of the manuscript.

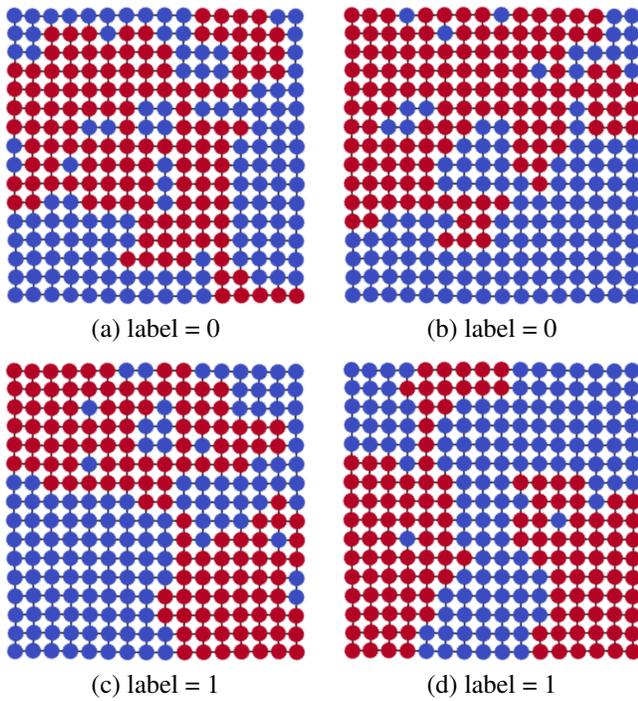
#### References

- [1] W.L. Hamilton, *Graph Representation Learning*, Morgan & Claypool Publishers, 2020.
- [2] M.M. Bronstein, J. Bruna, T. Cohen, and P. Veličković, “Geometric deep learning: Grids, groups, graphs, geodesics, and gauges,” *arXiv:2104.13478*, 2021.
- [3] J. Gilmer, S.S. Schoenholz, P.F. Riley, O. Vinyals, and G.E. Dahl, “Neural message passing for quantum chemistry,” in *Proc. of ICML*, 2017, PMLR, pp. 1263–1272.
- [4] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, “How powerful are graph neural networks?,” in *Proc. of ICLR*, 2019.
- [5] U. Alon and E. Yahav, “On the bottleneck of graph neural networks and its practical implications,” *Proc. of ICLR*, 2020.
- [6] Q. Li, Z. Han, and X.-M. Wu, “Deeper insights into graph convolutional networks for semi-supervised learning,” in *AAAI Conference on Artificial Intelligence*, 2018.
- [7] Y. Min, F. Wenkel, and G. Wolf, “Scattering GCN: Overcoming oversmoothness in graph convolutional networks,” in *Adv. in NeurIPS* 33, 2020, pp. 14498–14508.
- [8] H. Gao and S. Ji, “Graph U-nets,” in *Proc. of ICML*, 2019, PMLR, pp. 2083–2092.
- [9] M. Li, S. Chen, Y. Zhang, and I. Tsang, “Graph cross networks with vertex infomax pooling,” in *Adv. in NeurIPS* 33, 2020.
- [10] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar, “Multipole graph neural operator for parametric partial differential equations,” in *Adv. in NeurIPS* 33, 2020.
- [11] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, “Hierarchical graph representation learning with differentiable pooling,” in *Adv. in NeurIPS* 31, 2018.
- [12] F. Diehl, “Edge contraction pooling for graph neural networks,” *arXiv:1905.10990*, 2019.
- [13] C. Deng, Z. Zhao, Y. Wang, Z. Zhang, and Z. Feng, “GraphZoom: A multi-level spectral approach for accurate and scalable graph embedding,” in *Proc. of ICLR*, 2020.
- [14] V.D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, “Fast unfolding of communities in large networks,” *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, 2008.
- [15] T.N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” in *Proc. of ICLR*, 2017.
- [16] M. Schlichtkrull, T.N. Kipf, P. Bloem, R. Van Den Berg, I. Titov, and M. Welling, “Modeling relational data with graph convolutional networks,” in *European semantic web conference*, 2018.
- [17] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, “Graph attention networks,” in *Proc. of ICLR*, 2018.
- [18] S. Tang, B. Li, and H. Yu, “ChebNet: Efficient and stable constructions of deep neural networks with rectified power units using chebyshev approximations,” *arXiv:1911.05467*, 2019.
- [19] F. Errica, M. Podda, D. Bacciu, and A. Micheli, “A fair comparison of graph neural networks for graph classification,” in *Proc. of ICLR*, 2020.
- [20] Z. Yang, W. Cohen, and R. Salakhudinov, “Revisiting semi-supervised learning with graph embeddings,” in *Proc. of ICML*, 2016, vol. 48 of *PMLR*, pp. 40–48.
- [21] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra, “Beyond homophily in graph neural networks: Current limitations and effective designs,” in *Adv. in NeurIPS* 33, 2020.
- [22] Q. Huang, H. He, A. Singh, S.-N. Lim, and A.R. Benson, “Combining label propagation and simple models outperforms graph neural networks,” in *Proc. of ICLR*, 2021.
- [23] C. Morris, N.M. Kriege, F. Bause, K. Kersting, P. Mutzel, and M. Neumann, “TUDataset: A collection of bench-

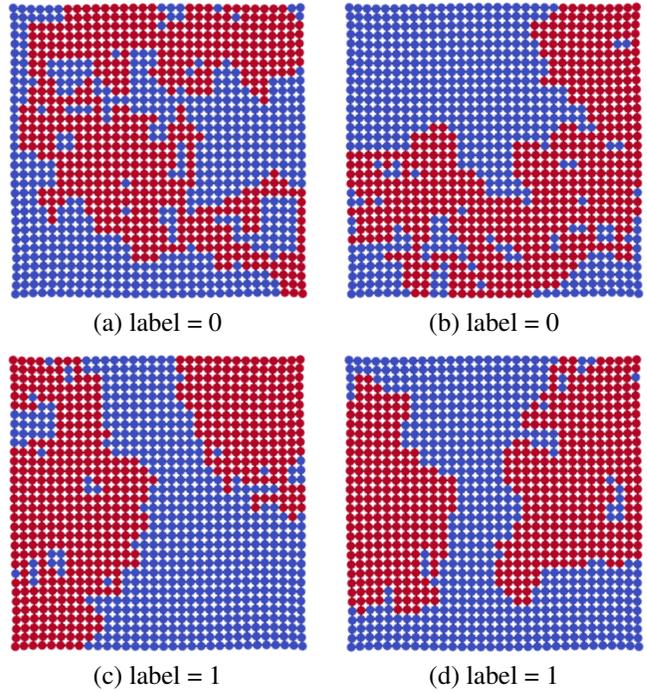
mark datasets for learning with graphs,” in *ICML 2020 GRL+ Workshop*, 2020.

- [24] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, “Open graph benchmark: Datasets for machine learning on graphs,” *Adv. in NeurIPS 33*, 2020.
- [25] K. Stachenfeld, J. Godwin, and P. Battaglia, “Graph networks with spectral message passing,” arXiv:2101.00079, 2020.
- [26] R.A. Rossi and N.K. Ahmed, “The network data repository with interactive graph analytics and visualization,” in *Proc. of AAAI*, 2015.

## Appendix



**Fig. A.1.** Two negative and two positive examples from 16x16 grid Color-connectivity dataset.



**Fig. A.2.** Two negative and two positive examples from 32x32 grid Color-connectivity dataset.

**Table A.1. Datasets summary.** The graph statistics are computed over all graphs in respective datasets. For model evaluation we used either the standard train/validation/test split as provided with the respective benchmark dataset and repeated the experiment 10 times with different random seeds (10x RS standard split); or we used 10-fold stratified cross-validation protocol (10-fold stratified CV). In transductive semi-supervised node classification with  $k$ -hop sanitized node pre-filtering we followed the same train/validation/test splitting procedure of [20]; train: 20 random pre-filtered nodes per each class, validation: 500 random nodes from remaining pre-filtered nodes, and test: 1000 random nodes from remaining pre-filtered nodes.

Dataset	# Graphs	(avg.) # Nodes	(avg.) # Edges	(Node) Features	# Classes	Evaluation	Metric
Cora	1	2,708	5,429	1,433	7	10x RS standard split	accuracy
CiteSeer	1	3,327	4,552	3,703	6	10x RS standard split	accuracy
PubMed	1	19,717	44,338	500	3	10x RS standard split	accuracy
COLLAB	5,000	74.49	2457.78	node degree	3	10-fold stratified CV	accuracy
IMDB-BINARY	1,000	19.77	96.53	node degree	2	10-fold stratified CV	accuracy
IMDB-MULTI	1,500	13	65.94	node degree	3	10-fold stratified CV	accuracy
D&D	1,178	284.32	715.66	89	2	10-fold stratified CV	accuracy
NCII	4,110	29.87	32.3	37	2	10-fold stratified CV	accuracy
ENZYMES	600	32.63	62.14	3	6	10-fold stratified CV	accuracy
PROTEINS	1,113	39.06	72.82	3	2	10-fold stratified CV	accuracy
ogbg-molpcba	437,929	26	28.1	9 node f. 3 edge f.	128 binary multilabel	10x RS standard split	avg. precision
ogbg-molhiv	41,127	25.5	27.5	9 node f. 3 edge f.	2	10x RS standard split	ROC-AUC
C-C 16x16 grid	15,000	256	480	1	2	10-fold stratified CV	accuracy
C-C 32x32 grid	15,000	1,024	1,984	1	2	10-fold stratified CV	accuracy
C-C Euroroad	15,000	1,174	1,417	1	2	10-fold stratified CV	accuracy
C-C Minnesota	6,000	2,642	3,304	1	2	10-fold stratified CV	accuracy