

Graph Neural Networks (GNNs)

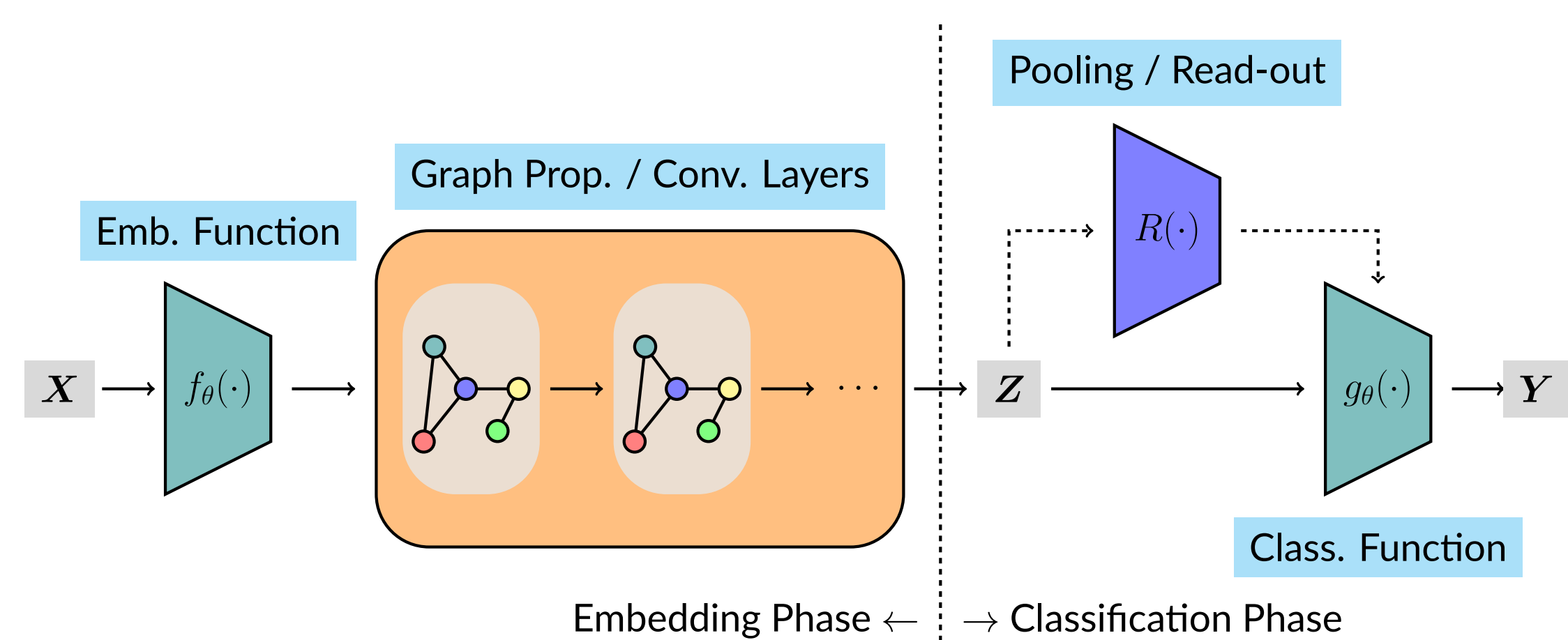


Figure 1. A typical structure of GNNs.

- Embedding function $f(\mathbf{X})$: map features $\mathbf{X} \in \mathbb{R}^{n \times F}$ to initial embeddings $\mathbf{Z}^{(0)} \in \mathbb{R}^{n \times d}$;
- Graph propagation / convolution $\text{GP}(\mathbf{Z}^{(k)}; \mathbf{A})$: propagate embedding K times;
- Read-out / pooling $\text{R}(\mathbf{Z}^{(K)}; \mathbf{A})$ (optional): pooling for graph-level tasks;
- Classification function $g(\cdot)$: final classification to generate predictions \mathbf{Y} .

Expressive Power & Universality of GNNs

- Spectral GNNs: designing **universal filters** (e.g., ChebNet, GPRGNN, BernNet, ...).
- Spatial GNNs: designing GNNs bounded by k -WL tests (e.g., GIN, ...).

There are works exploring relations between GNNs and geometric objects (e.g., curvature, cellular sheaves) and physical concepts (e.g., oscillators).

Q: Can we define universality of spatial GNNs from a geometric perspective?

Preliminaries

Definition (Equivalent). For a given graph $G = (V, E)$, two node embedding matrices $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)} \in \mathbb{R}^{n \times d}$ are equivalent if for all $(i, j) \in E$, $\|\mathbf{Z}_i^{(1)} - \mathbf{Z}_j^{(1)}\|_2 = \|\mathbf{Z}_i^{(2)} - \mathbf{Z}_j^{(2)}\|_2$ holds.

Definition (Congruent). For a given graph $G = (V, E)$, two node embedding matrices $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)} \in \mathbb{R}^{n \times d}$ are congruent if for all $i, j \in V$, $\|\mathbf{Z}_i^{(1)} - \mathbf{Z}_j^{(1)}\|_2 = \|\mathbf{Z}_i^{(2)} - \mathbf{Z}_j^{(2)}\|_2$ holds.

Definition (Globally Rigid). For a given graph $G = (V, E)$, an embedding matrix \mathbf{Z} is globally rigid if all its equivalent embedding matrices \mathbf{Z}' are also congruent to \mathbf{Z} .

Definition (Rigid). For a given graph $G = (V, E)$, an embedding matrix \mathbf{Z} is rigid if all equivalent embeddings that can be obtained by **continuous motion** from \mathbf{Z} are congruent to \mathbf{Z} .

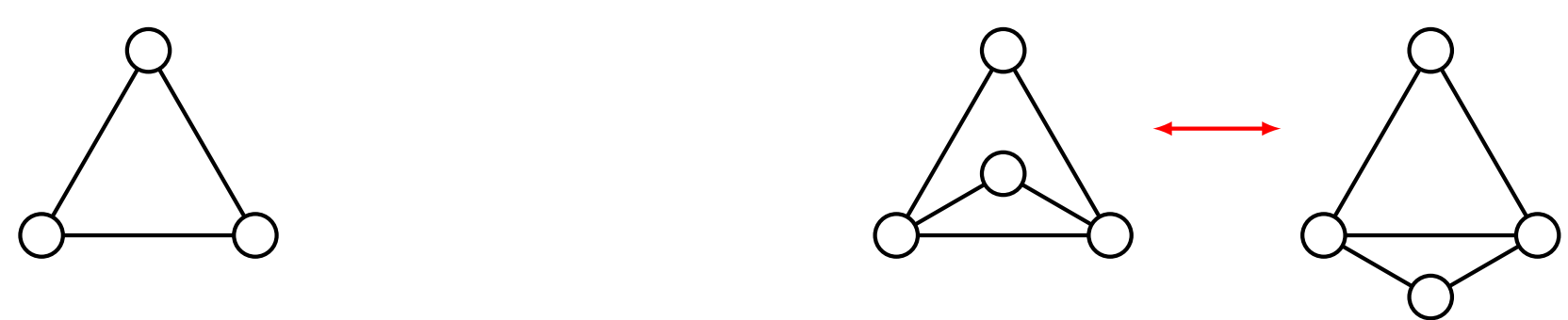
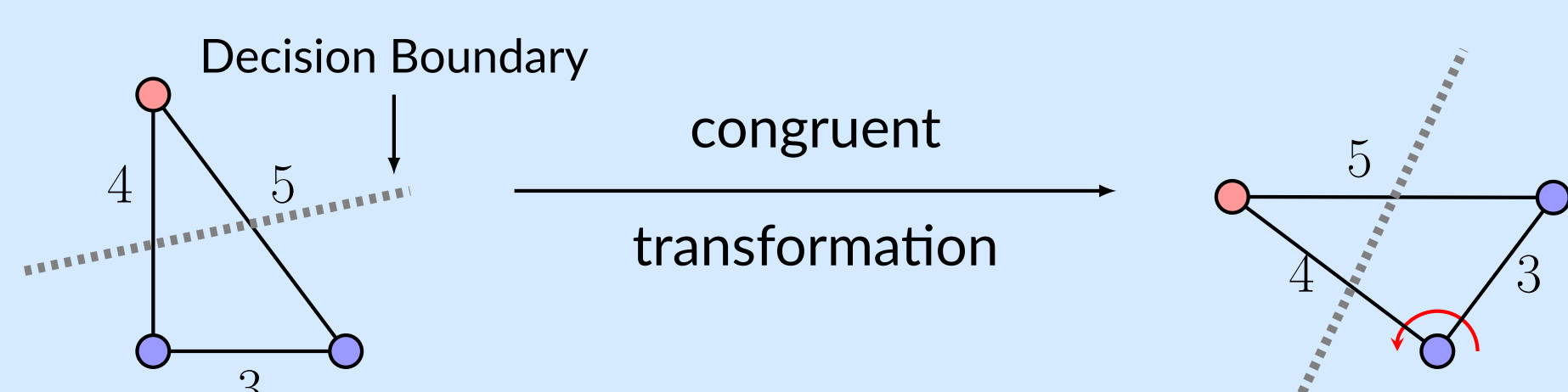


Figure 2. Left: A globally rigid graph in \mathbb{R}^2 ; Right: A rigid but not globally rigid graph in \mathbb{R}^2 , since it has an equivalent but not congruent embedding.

Definition (Metric Matrix). The metric matrix of an embedding matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$ is defined as $\mathbf{M}_{\mathbf{Z}} = (\|\mathbf{Z}_i - \mathbf{Z}_j\|_2)_{ij}$. We also define the mapping from an embedding matrix \mathbf{Z} to its metric matrix $\mathbf{M}_{\mathbf{Z}}$ as $\mathbf{M}_{\mathbf{Z}} = \mathbf{M}(\mathbf{Z})$.

Defining Spatial-Universality

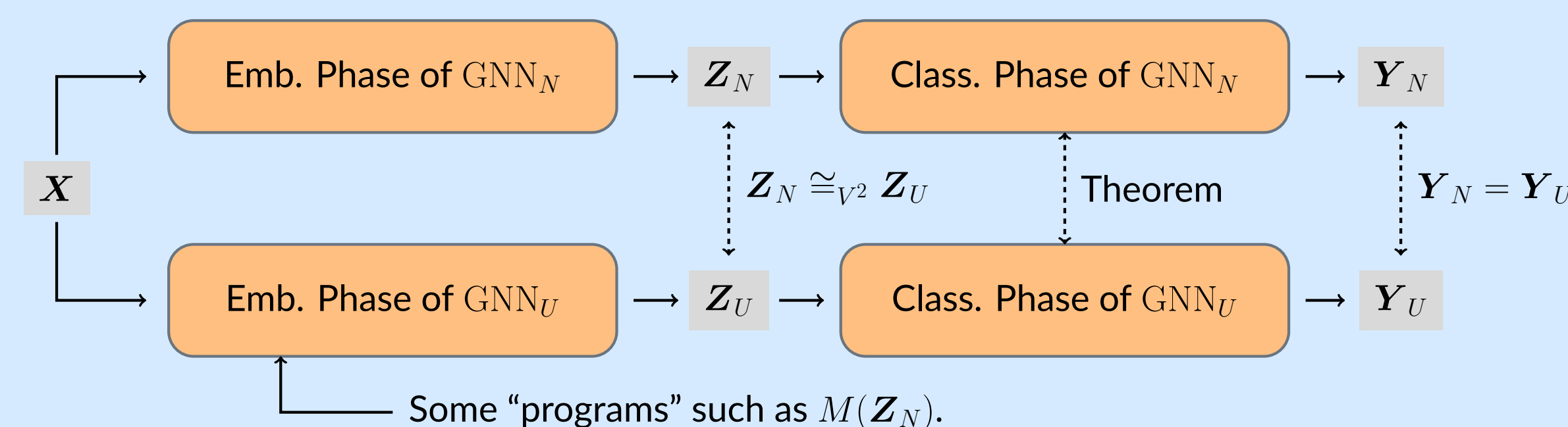
Observation:



Defining Spatial-Universality (Cont'd)

Theorem (MLPs Are Congruent-Insensitive). Given two congruent embedding matrices \mathbf{Z}_1 and \mathbf{Z}_2 , for any MLP_M (with biases), there always exists another MLP_N (also with biases) such that $\text{MLP}_M(\mathbf{Z}_1) = \text{MLP}_N(\mathbf{Z}_2)$.

Idea:



Spatial-Universal: It can arrange nodes with a given metric matrix!

This idea is closely related to the **Distance Geometry Problem (DGP)**.

Distance Geometry Problem (DGP)

Given a positive integer d , a graph $G = (V, E)$, and a symmetric non-negative matrix \mathbf{M} , decide whether there exists an embedding matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$, such that

$$\forall (i, j) \in E, \|\mathbf{Z}_i - \mathbf{Z}_j\| = M_{ij}.$$

Optimization Objective

About the Optimization Objective

- **Full** metric matrix: $O(n^2) \Rightarrow$ **partial** metric matrix on edges: $O(m)$.
- For **globally rigid** graphs, partial metric matrix is enough to determine the “shape” of the embedding, this modification does not weaken the expressive power.
- However, solving the DGP is **NP-Hard**. We can not directly arrange the nodes, thus we introduce an **error-tolerant objective**:

$$E_p(\mathbf{Z}; \mathbf{M}, E) = \frac{1}{2} \|\mathbf{A} \odot (\mathbf{M}(\mathbf{Z}) - \mathbf{M})\|_F^2 = \sum_{(i,j) \in E} \frac{1}{2} (\|\mathbf{Z}_i - \mathbf{Z}_j\|_2 - M_{ij})^2.$$

- This function is closely related to the raw **Stress function** σ_r in the **Multidimensional Scaling (MDS)** problem and the potential energy of **spring networks**.

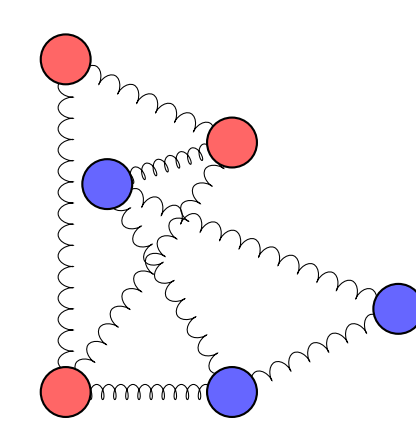


Figure 3. A spring network.

- To align with other representative GNNs, we modify E_p and add a regularization term to get the final objective:

$$\begin{aligned} \mathcal{L}(\mathbf{Z}; \mathbf{Z}^{(0)}, \mathbf{M}, E) &= (1 - \alpha) \tilde{E}_p(\mathbf{Z}; \mathbf{M}, E) + \alpha \|\mathbf{Z} - \mathbf{Z}^{(0)}\|_F^2 \\ &= (1 - \alpha) E_p(\mathbf{D}^{1/2} \mathbf{Z}; \mathbf{M}, E) + \alpha \|\mathbf{Z} - \mathbf{Z}^{(0)}\|_F^2. \end{aligned}$$

About the Metric Matrix

- For scenarios with prior knowledge about distances between nodes (e.g., molecular conformation generation, or graph drawing), directly use them; for other scenarios, learn a metric matrix.
- Our idea is to increase the distances between dissimilar nodes and reduce the distances between similar nodes:
 1. Introduce edge attention $\alpha_{ij} \in [-1, 1]$, when $\alpha_{ij} \rightarrow 1 \Leftrightarrow i, j$ tend to belong to the same class, and when $\alpha_{ij} \rightarrow -1 \Leftrightarrow i, j$ tend to belong to different classes;
 2. Map the initial embedding matrix $\mathbf{Z}^{(0)}$ (defined later) to a hidden matrix \mathbf{H} ;
 3. Use attention mechanisms, such as $\alpha_{ij} = \tanh(\mathbf{a}^\top [\mathbf{H}_i^\top; \|\mathbf{H}_j^\top\|])$ or $\alpha_{ij} = \tanh(\mathbf{H}_i \cdot \mathbf{W} \mathbf{H}_j^\top)$ to learn the edge attention;
 4. Then we can set $M_{ij} = \frac{1 - \alpha_{ij}}{1 + \alpha_{ij} + \varepsilon} \|\mathbf{Z}_i^{(0)} - \mathbf{Z}_j^{(0)}\|$, where ε is a small positive number.

Framework

The Embedding Function

- A linear layer $f(\mathbf{X}) = \mathbf{X} \mathbf{W} + \mathbf{1} \mathbf{b}^\top$ in linear GNNs, or
- A two-layer MLP $f(\mathbf{X}) = \sigma(\sigma(\mathbf{X} \mathbf{W}_1 + \mathbf{1} \mathbf{b}_1^\top) \mathbf{W}_2 + \mathbf{1} \mathbf{b}_2^\top)$ in spectral GNNs.

Propagation

- Our goal is to design a propagation method that minimizes the objective.
- It's non-convex. Following related works, we employ the **stationary point iteration** method.
- By computing the gradient, setting it to zero, rearranging the terms, rewriting it as an iteration form, substituting $1 - \alpha$ with β to allow more flexibility, it leads to:

$$\mathbf{Z}^{(k+1)} = (1 - \alpha) \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \mathbf{Z}^{(k)} + \beta \mathbf{D}^{-1/2} \mathbf{L}_H \mathbf{D}^{-1/2} \mathbf{Z}^{(k)} + \alpha \mathbf{Z}^{(0)},$$

where $\mathbf{H} = \mathbf{A} \odot \mathbf{M} \odot \mathbf{M}(\mathbf{D}^{-1/2} \mathbf{Z})^{\odot -1}$, and $\mathbf{L}_H = \text{diag}(\mathbf{H} \mathbf{1}) - \mathbf{H}$.

- And we have the message-passing form:

$$\mathbf{Z}_i^{(k+1)} = (1 - \alpha) \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{Z}_j^{(k)}}{\sqrt{d_i d_j}} + \beta \sum_{j \in \mathcal{N}(i)} \frac{M_{ij} (\mathbf{Z}_i^{(k)} - \mathbf{Z}_j^{(k)})}{\sqrt{d_i d_j} \|\mathbf{Z}_i^{(k)} / \sqrt{d_i} - \mathbf{Z}_j^{(k)} / \sqrt{d_j}\|_2} + \alpha \mathbf{Z}_i^{(0)}.$$

Optional Linear and Non-linear Transformations

We may incorporate linear and non-linear transformations after each propagation step. In our experiments without pre-designed metric matrices, such as node classification, we utilize three designs from the GCNII model: a linear transformation, the identity mapping, and a non-linear transformation (ReLU).

The Classification Function

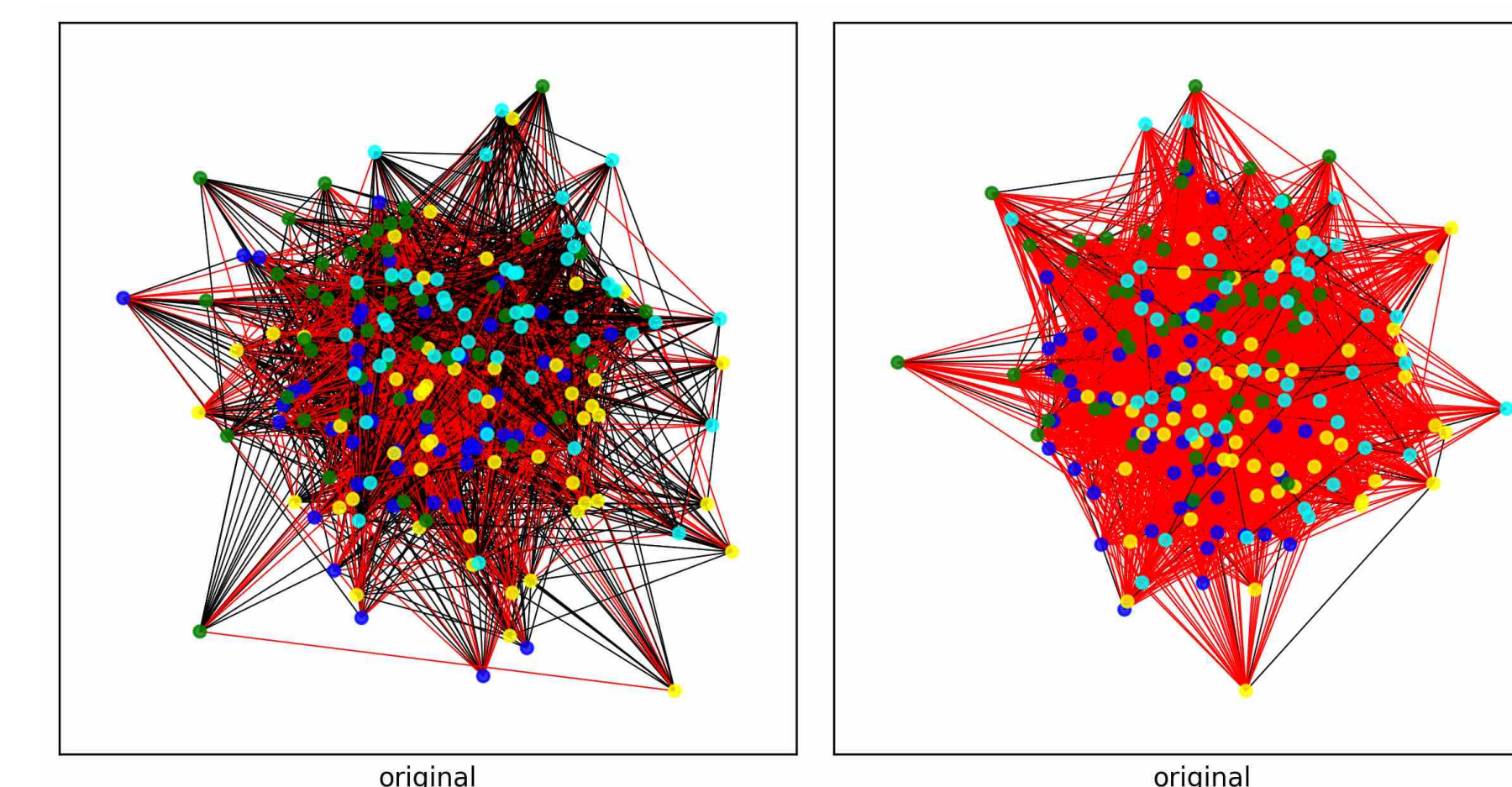
We use a linear layer $g(\mathbf{Z}^{(K)}) = \mathbf{Z}^{(K)} \mathbf{W} + \mathbf{1} \mathbf{b}^\top$ as the final classification function.

Experiments

We have done the “Arranging Nodes with Given Metric Matrices” experiments on synthetic graphs, supervised node classification and graph regression experiments on real-world graphs.

Arranging Nodes with Given Metric Matrices

- We generate two **Stochastic Block Model (SBM)** graphs, one homophilic and one heterophilic, consisting of four blocks with 50 nodes in each block.
- The node features are sampled from two 2-dimensional Gaussian distributions.



- If i and j are in the same class, we set $M_{ij} = 0$; otherwise, we set $M_{ij} = 5$.
- We pass the node features through 8 MGNN layers, with $\alpha = 0.05$ and $\beta = 0.5$.
- For visualization results, please refer to the smaller posters below or our paper.

Supervised Node Classification and Graph Regression

- Our MGNN model performs well, and the results are promising.
- For experiment details and results, please refer to the smaller posters below or our paper.

Experiment I: Arranging Nodes with Given Metric Matrices

- We pass the node features through 8 MGNN layers, with $\alpha = 0.05$ and $\beta = 0.5$. After each round of propagation, we visualize the results.

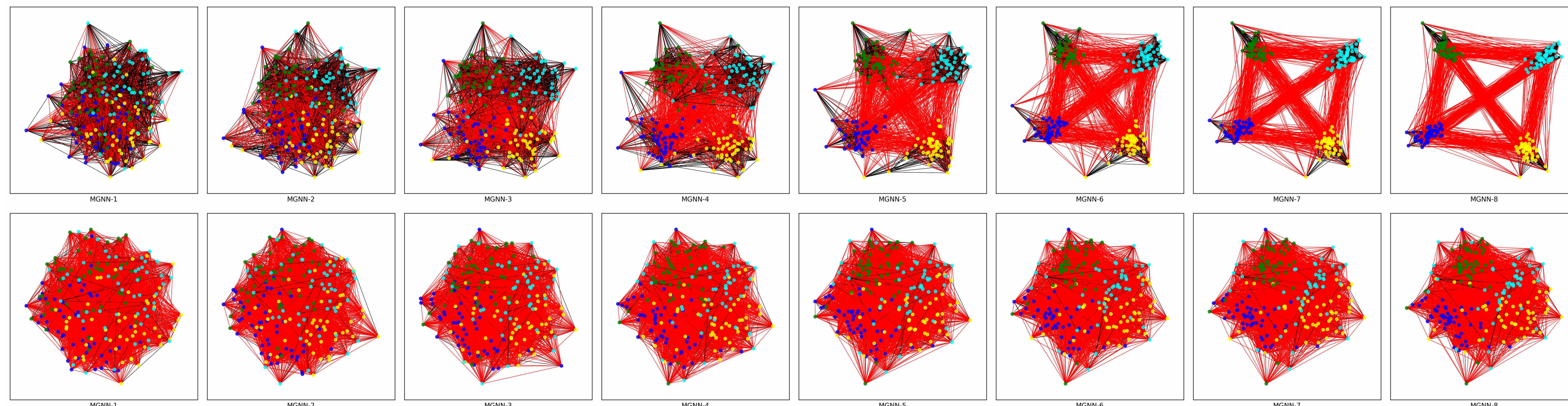


Figure 1. Visualization results of MGNN propagation layers.

- Based on the results, we can confirm that our MGNN model aims to separate the four blocks.
- The presence of homophilic edges adds convexity to our objective function, which is easier to optimize.
- Additionally, we provide visualization results obtained from the SGC layers and APPNP layers.

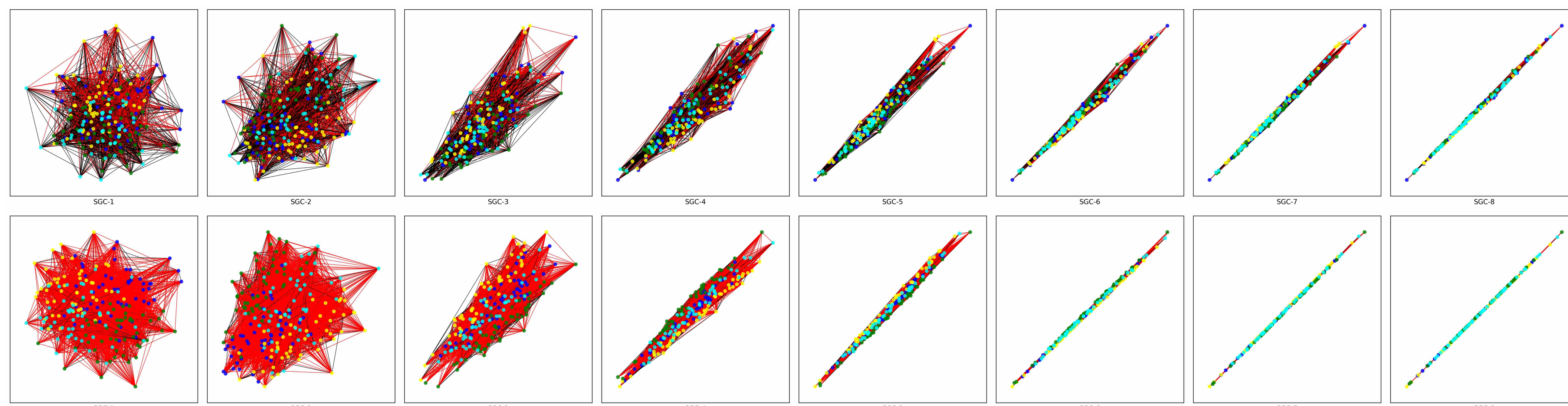


Figure 2. Visualization results of SGC propagation layers.

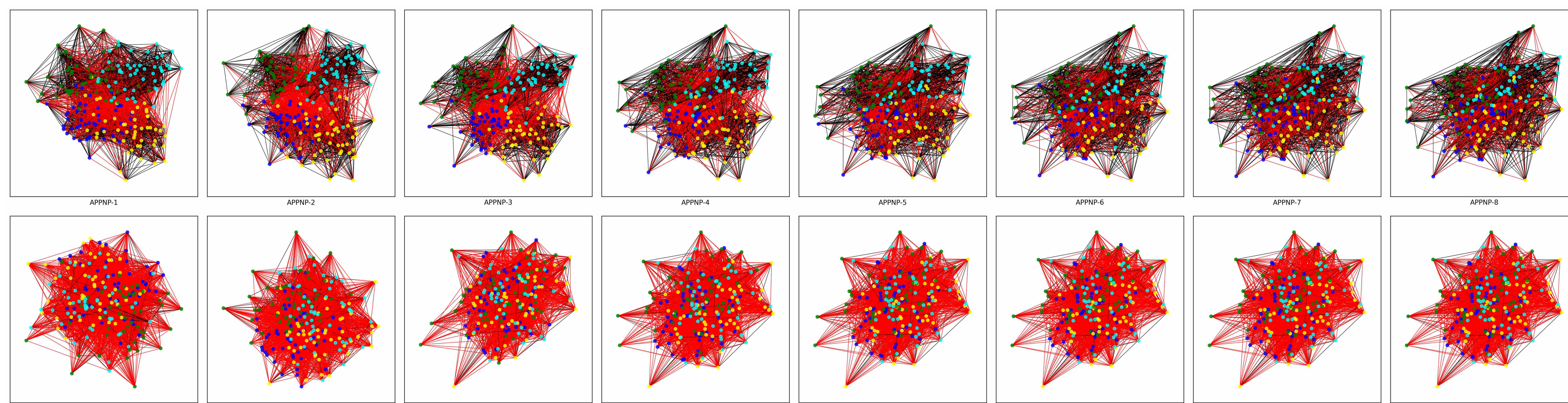


Figure 3. Visualization results of APPNP propagation layers.

Experiment II: Supervised Node Classification

- For all datasets except ogbn-arxiv, we generate 10 random splits with train/valid/test ratio of 60%/20%/20%. For the ogbn-arxiv dataset, we generate 10 random splits using the `get_idx_split()` function provided by the official package.
- We then tune the hyper-parameters 100 times for each model on each dataset. (Please refer to our paper for the detailed ranges of hyperparameters.)
- We train each model for a maximum of 1500 epochs, with early stopping set to 100, on each dataset and report the results after hyper-tuning.

Table 1. The results (accuracy with standard deviation) of the supervised node classification experiments. Boldface results indicate the best model on each dataset, and underlined results are second best models.

	Non-Graph		Spectral			Spatial (Non-Spectral)					MGNN	
	Linear	MLP	GCN	SGC	APPNP	PointNet	GAT	GIN	GCNII	pGNN		LINKX
Cora	76.09±1.55	76.11±1.58	88.25±1.09	88.34±1.41	89.30±1.45	84.43±1.94	88.71±0.89	85.76±1.18	88.52±1.40	88.78±1.20	82.82±1.96	89.04±1.20
CiteSeer	71.11±1.81	73.25±1.16	77.10±1.12	<u>77.47±1.42</u>	76.80±1.10	72.83±1.38	76.44±1.19	72.83±1.47	77.05±0.78	77.28±0.88	71.79±1.55	78.08±1.50
PubMed	87.06±0.67	88.21±0.46	89.07±0.42	<u>87.59±0.54</u>	89.75±0.47	89.18±0.38	87.77±0.67	87.15±0.54	<u>90.24±0.55</u>	89.79±0.38	86.77±0.59	90.37±0.47
CoraFull	60.55±0.76	61.61±0.56	71.86±0.75	71.86±0.70	71.88±0.75	63.32±1.02	70.65±0.86	67.11±0.46	<u>71.09±0.72</u>	<u>72.36±0.58</u>	63.57±0.55	72.42±0.69
CS	94.51±0.32	94.89±0.23	93.77±0.37	94.10±0.40	95.91±0.23	93.13±0.42	93.25±0.31	91.81±0.34	<u>95.98±0.22</u>	95.83±0.23	95.06±0.24	96.01±0.16
Physics	95.92±0.15	96.01±0.26	96.46±0.25	OOM	<u>97.14±0.21</u>	96.37±0.27	96.47±0.21	94.66±0.24	97.10±0.21	96.93±0.10	96.87±0.16	97.15±0.15
Cornell	<u>78.65±2.82</u>	75.95±4.43	50.27±6.86	53.51±4.32	77.84±6.49	71.08±5.93	50.27±7.66	49.46±7.46	76.76±7.76	77.30±8.22	71.08±12.09	81.89±6.29
Texas	81.35±4.90	83.78±7.55	61.08±8.65	56.22±6.37	86.76±1.46	82.16±6.30	61.08±5.57	64.05±4.20	<u>88.65±4.95</u>	85.14±5.16	84.32±5.24	90.00±2.72
Wisconsin	86.20±3.63	88.80±2.40	55.80±5.83	58.00±4.29	86.00±3.10	81.60±3.98	56.40±5.99	57.20±6.82	88.20±4.14	85.40±3.58	82.00±3.10	88.40±3.44
Chameleon	50.77±2.07	50.75±2.13	<u>69.21±2.08</u>	67.12±2.17	67.85±2.65	63.76±2.52	66.97±2.45	46.73±13.51	67.56±1.18	69.19±1.57	67.47±1.62	72.37±2.25
Squirrel	35.76±1.05	35.79±1.65	55.43±2.05	52.18±1.49	54.60±1.88	47.39±8.31	<u>55.68±2.81</u>	20.96±2.08	53.88±2.77	51.61±1.28	57.86±1.17	54.45±1.85
Actor	36.16±0.75	37.67±1.60	30.42±1.58	30.14±1.18	36.98±1.28	36.41±1.23	29.22±0.94	26.09±1.75	<u>37.75±1.23</u>	36.47±1.07	35.00±2.11	38.36±1.46
WikiCS	78.74±0.57	79.85±0.69	84.02±0.61	83.47±0.83	84.88±0.55	84.09±0.86	83.82±0.73	66.08±22.77	<u>85.09±0.71</u>	84.41±0.46	84.13±0.56	85.09±0.59
ogbn-arxiv	52.80±0.19	53.81±0.23	70.48±0.18	68.77±0.06	70.50±0.11	69.89±0.59	<u>71.04±0.27</u>	66.60±0.53	71.48±0.21	68.50±0.17	59.28±2.22	70.73±0.17
Average Rank	9.21	7.86	6.79	7.69	3.50	7.86	7.57	10.50	3.50	4.36	7.29	1.57

Experiment III: Graph Regression

- We utilized the publicly available train/validation/test split of the ZINC-subset dataset.
- For all models, we set the number of hidden units to 64. We limit the training of each model to a maximum of 500 epochs. A batch size of 512 is utilized, and early stopping is implemented after 50 epochs on each dataset.
- We conduct hyperparameter tuning for each model on each dataset and repeat this process 100 times. (Please refer to our paper for the detailed ranges of hyperparameters.)
- To enhance reliability, we repeat the whole process five times and calculate the mean MAE (Mean Absolute Error) and standard deviation for each model.

Table 2. The results of the graph regression experiments.

GCN	GAT	GIN	PointNet	MGNN
0.6143±0.0200	0.6458±0.0647	0.4410±0.0065	0.5293±0.0138	<u>0.4751±0.0112</u>