



**ICML**  
International Conference  
On Machine Learning

# Towards Deep Attention in GNNs: Problems and Remedies



**Soo Yong Lee**  
KAIST



**Fanchen Bu**  
KAIST



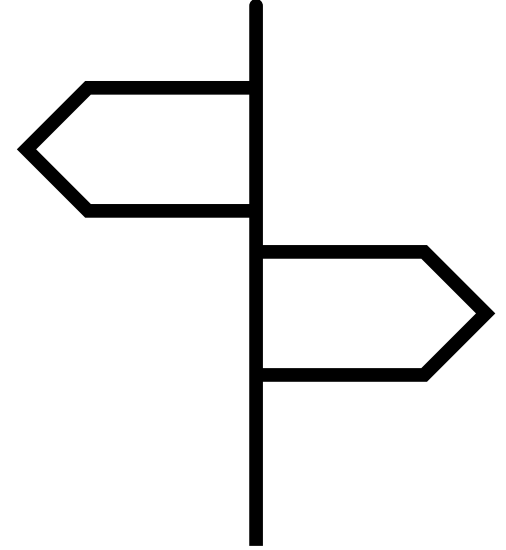
**Jaemin Yoo**  
CMU



**Kijung Shin**  
KAIST

# Contents

- Sec. 1: Introduction
- Sec. 2: Analysis of Graph Attention
- Sec. 3: Proposed Method : AERO-GNN
- Sec. 4: Experiments and Empirical Evaluation
- Sec. 5: Discussion



# Graphs

- **What are Graphs?**
  - Graphs are relational data
  - Consists of nodes and edges
- **Graphs are everywhere!**
  - Can represent a wide range of real-world networks



**Web Networks**  
Node = Webpage  
Edge = Hyperlinks



**Social Networks**  
Node = User  
Edge = Follow



**Transportation Networks**  
Node = Region  
Edge = Road Connection

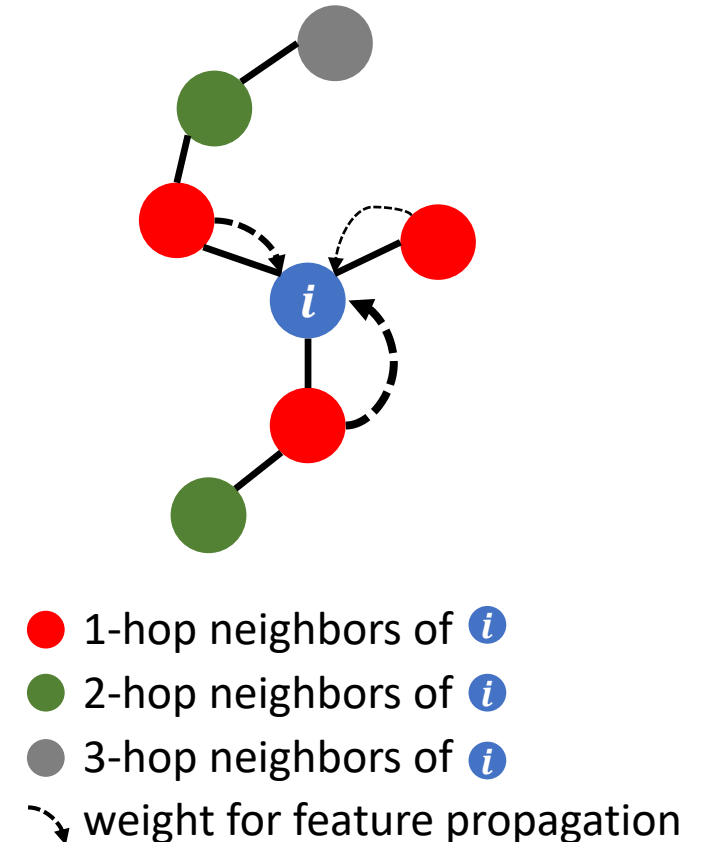


**Co-Purchase Networks**  
Node = Product  
Edge = Often Co-Purchased

# Graph Neural Networks (GNNs)

- **Graph Neural Networks (GNNs)**
  - Can solve various graph-related tasks
  - Learn graph representation
- **To enhance its expressiveness:**
  - Graph Attention
    - Learns the weight for feature propagation
  - Deep GNN
    - Increases receptive fields
    - Stacks non-linearity

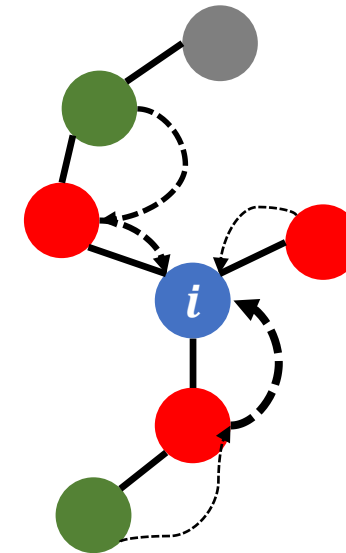
1-Layer Graph Attention:  
Receptive Field of Node  $i$



# Graph Neural Networks (GNNs)

- **Graph Neural Networks (GNNs)**
  - Can solve various graph-related tasks
  - Learn graph representation
- **To enhance its expressiveness:**
  - Graph Attention
    - Learns the weight for feature propagation
  - Deep GNN
    - Increases receptive fields
    - Stacks non-linearity

**2-Layer** Graph Attention:  
Receptive Field of Node  $i$



- 1-hop neighbors of  $i$
- 2-hop neighbors of  $i$
- 3-hop neighbors of  $i$
- ↪ weight for feature propagation

# Goal of the Present Study



## Question of Interest

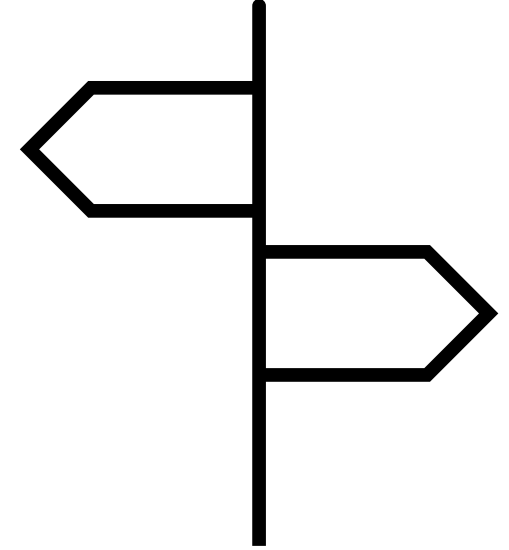
Can existing graph attention **remain expressive** over deep layers?

How to design an **expressive deep graph attention**?

Can it solve **node classification** problem?

# Contents

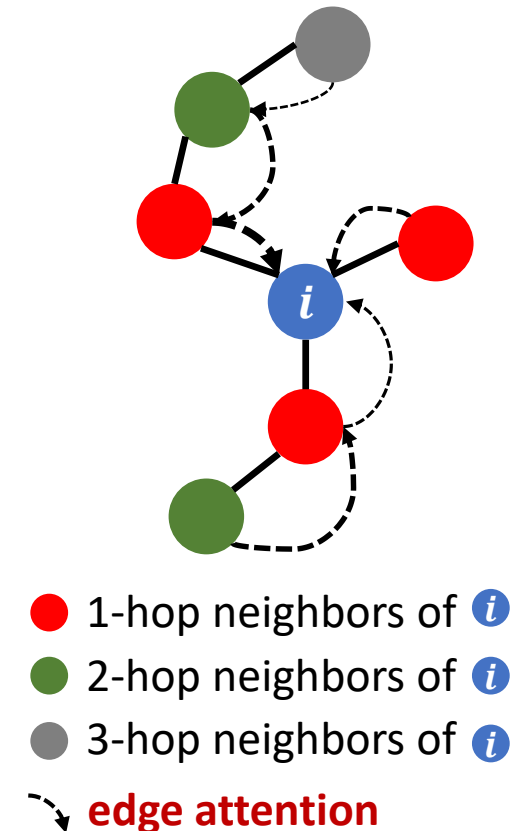
- Sec. 1: Introduction
- Sec. 2: Analysis of Graph Attention
- Sec. 3: Proposed Method : AERO-GNN
- Sec. 4: Experiments and Empirical Evaluation
- Sec. 5: Discussion



# Graph Attention for GNNs

- **Edge Attention**  $A^{(k)}$ 
  - Intuition: learns importance *within* each hop
  - Models: GAT[1], FAGCN[2]
- **Hop Attention**  $\Gamma^{(k)}$ 
  - Intuition: learns importance *of* each hop
  - Models: GPRGNN[3], DAGNN[4]

Illustration of Hop Attention



[1] Velickovic,P., Cucurull,G., Casanova,A., Romero,A., Lio,P., and Bengio,Y. Graph attention networks. In ICLR, 2018.

[2] Bo,D., Wang,X., Shi,C., and Shen,H. Beyond low frequency information in graph convolutional networks. In AAAI, 2021.

[3] Liu,M., Gao,H., and Ji,S. Towards deeper graph neural networks. In KDD, 2020.

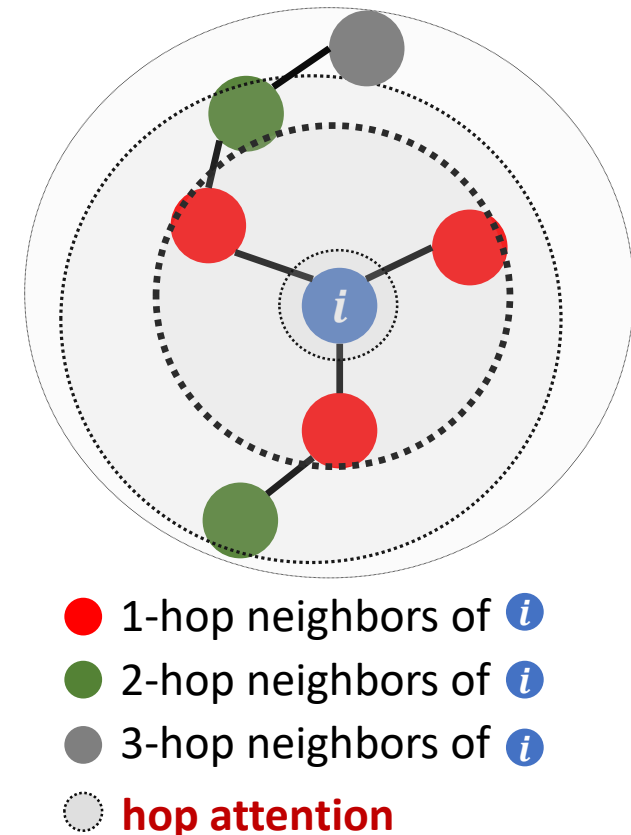
[4] Chien,E., Peng,J., Li,P., and Milenkovic,O. Adaptive universal generalized pagerank graph neural network. In ICLR, 2021.



# Graph Attention for GNNs

- **Edge Attention  $A^{(k)}$** 
  - Intuition: learns importance *within* each hop
  - Models: GAT[1], FAGCN[2]
- **Hop Attention  $\Gamma^{(k)}$** 
  - Intuition: learns importance *of* each hop
  - Models: GPRGNN[3], DAGNN[4]

Illustration of Hop Attention



[1] Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. Graph attention networks. In ICLR, 2018.

[2] Bo, D., Wang, X., Shi, C., and Shen, H. Beyond low frequency information in graph convolutional networks. In AAAI, 2021.

[3] Liu, M., Gao, H., and Ji, S. Towards deeper graph neural networks. In KDD, 2020.

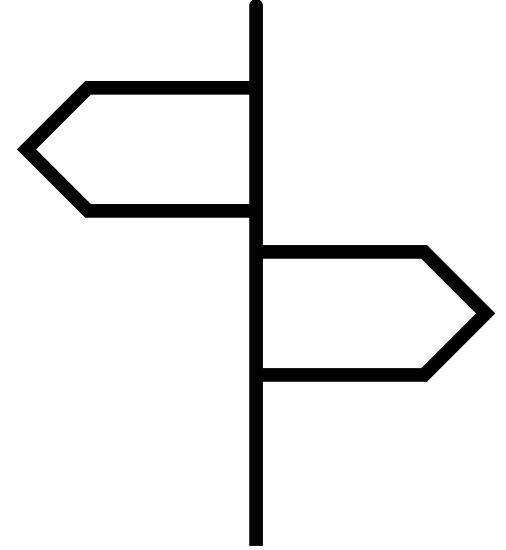
[4] Chien, E., Peng, J., Li, P., and Milenkovic, O. Adaptive universal generalized pagerank graph neural network. In ICLR, 2021.

# Theoretical Limitations to Deep Graph Attention

- **All Graph Attention Models Suffer From Two Problems**
  - P1: Vulnerability of Node Feature Over-Smoothing
    - (Informal) The **attention coefficients become identical** for over-smoothed node features
  - P2: Smooth Cumulative Attention
    - (Informal) **Cumulative attention vectors become identical** for all nodes at very deep layer
- ***Both problems are critically contrary to the goal of attention***

# Contents

- Sec. 1: Introduction
- Sec. 2: Analysis of Graph Attention
- **Sec. 3: Proposed Method : AERO-GNN**
- Sec. 4: Experiments and Empirical Evaluation
- Sec. 5: Discussion



# AERO-GNN : Overview

- We propose Attentive Deep Propagation GNN (**AERO-GNN**)
- **Model Overview**
  - At every propagation layer  $k$ , AERO-GNN learns  $\mathbf{A}^{(k)}$  and  $\mathbf{\Gamma}^{(k)}$

$$H^{(k)} = \begin{cases} \text{MLP}(X), & \text{if } k = 0, \\ \underline{\mathbf{A}}^{(k)} H^{(k-1)}, & \text{if } 1 \leq k \leq k_{max}, \end{cases}$$

$$Z^{(k)} = \sum_{l=0}^k \underline{\mathbf{\Gamma}}^{(l)} H^{(l)}, \forall 1 \leq k \leq k_{max},$$

$$Z^* = \sigma(Z^{(k_{max})})W^*,$$

# AERO-GNN : Attention Functions

- **Design Question :**

- How do we design an expressive deep graph attention?

- **Key Properties :**

- Key 1. Stacking **non-linearity**
- Key 2. Learn both  $A^{(k)}$  and  $\Gamma^{(k)}$  (edge and hop attention)
- Key 3. Use features from the previous layers  $Z$
- Key 4. Use **negative** attention
- Key 5. Have **node-adaptive** hop attention  $\Gamma^{(k)}$

# AERO-GNN : Attention Functions

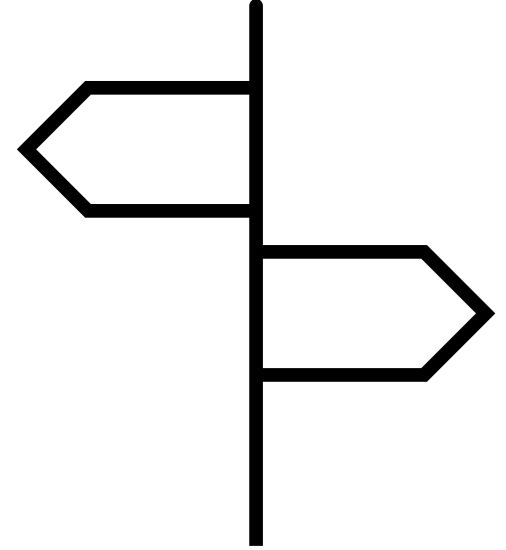
## ■ Bottom Line :

- *Attention functions of AERO-GNN is flexible and expressive!*
- *They allow AERO-GNN to mitigate problems of deep graph attention.*
  - Vulnerability to Over-Smoothing & Smooth Cumulative Attention

Properties of Attention Functions					
	Stacking Non-Linearity	Edge & Hop	Z as Input	Negative Attention	Node-Adaptive
GATv2	○	✗	✗	✗	✗
FAGCN	○	✗	○	○	✗
GPRGNN	✗	✗	✗	○	✗
DAGNN	✗	✗	✗	✗	○
AERO-GNN	○	○	○	○	○

# Contents

- Sec. 1: Introduction
- Sec. 2: Analysis of Graph Attention
- Sec. 3: Proposed Method : AERO-GNN
- Sec. 4: Experiments and Empirical Evaluation
- Sec. 5: Discussion



# Performance (Mean $\pm$ Std, 100 trials)

- **AERO-GNN achieves the best overall performance (See high A.R.)!**

Table 3: Node Classification Performance on Real-World Graphs

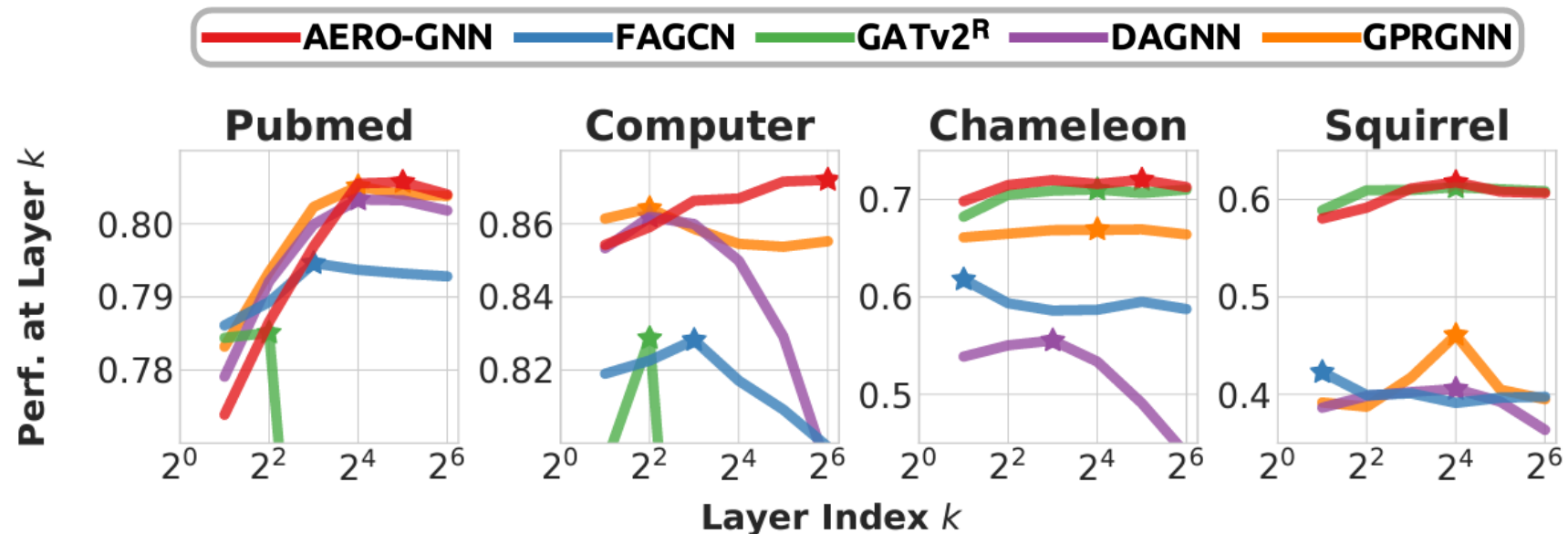
Dataset	Chameleon	Squirrel	Actor	Texas	Cornell	Wisconsin	Computer	Photo	Wiki-CS	Pubmed	Citeseer	Cora	A.R.
<b>Homophily</b>	0.04	0.03	0.01	0.00	0.02	0.05	0.70	0.77	0.57	0.66	0.63	0.77	
<b>GCN</b>	67.97 $\pm$ 2.5	53.33 $\pm$ 1.3	30.57 $\pm$ 0.7	65.65 $\pm$ 4.8	58.41 $\pm$ 3.3	62.02 $\pm$ 5.9	81.27 $\pm$ 1.4	90.24 $\pm$ 1.3	79.08 $\pm$ 0.5	79.54 $\pm$ 0.4	72.50 $\pm$ 0.5	83.15 $\pm$ 0.5	9.1
<b>APPNP</b>	53.04 $\pm$ 2.2	40.37 $\pm$ 1.5	35.49 $\pm$ 1.0	79.89 $\pm$ 4.2	80.16 $\pm$ 5.9	84.24 $\pm$ 4.6	81.27 $\pm$ 1.4	91.12 $\pm$ 1.2	79.05 $\pm$ 0.5	79.90 $\pm$ 0.3	73.06 $\pm$ 0.3	83.60 $\pm$ 1.3	7.8
<b>GCN-II</b>	60.38 $\pm$ 1.9	48.76 $\pm$ 2.4	35.77 $\pm$ 1.0	78.59 $\pm$ 6.6	78.84 $\pm$ 6.6	83.20 $\pm$ 4.7	84.24 $\pm$ 1.2	91.81 $\pm$ 0.9	79.28 $\pm$ 0.6	80.14 $\pm$ 0.6	<b>73.20 <math>\pm</math> 0.8</b>	<b>85.33 <math>\pm</math> 0.5</b>	5.5
<b>A-DGN</b>	69.63 $\pm$ 2.0	57.77 $\pm$ 1.9	36.41 $\pm$ 1.0	82.22 $\pm$ 4.8	<b>83.14 <math>\pm</math> 6.7</b>	<b>85.84 <math>\pm</math> 4.0</b>	83.70 $\pm$ 1.5	90.53 $\pm$ 1.3	79.11 $\pm$ 0.6	78.68 $\pm$ 0.6	70.16 $\pm$ 0.9	79.84 $\pm$ 0.9	6.4
<b>GAT</b>	68.01 $\pm$ 2.5	54.49 $\pm$ 1.7	30.36 $\pm$ 0.9	60.46 $\pm$ 6.2	58.22 $\pm$ 3.7	63.59 $\pm$ 6.1	84.46 $\pm$ 1.3	89.88 $\pm$ 1.1	79.44 $\pm$ 0.5	78.94 $\pm$ 0.4	71.89 $\pm$ 0.6	83.78 $\pm$ 0.5	8.5
<b>GATv2</b>	69.06 $\pm$ 2.2	57.67 $\pm$ 2.4	30.27 $\pm$ 0.8	60.32 $\pm$ 7.0	58.35 $\pm$ 3.8	61.94 $\pm$ 4.7	84.19 $\pm$ 1.2	89.87 $\pm$ 1.2	79.64 $\pm$ 0.5	79.12 $\pm$ 0.3	71.15 $\pm$ 1.2	83.88 $\pm$ 0.6	8.9
<b>GATv2<sup>R</sup></b>	<b>70.88 <math>\pm</math> 1.9</b>	<b>61.23 <math>\pm</math> 1.5</b>	33.73 $\pm$ 0.9	60.68 $\pm$ 6.6	57.32 $\pm$ 4.5	60.61 $\pm$ 5.1	81.73 $\pm$ 2.2	88.71 $\pm$ 1.7	79.75 $\pm$ 0.6	78.28 $\pm$ 0.4	71.00 $\pm$ 0.8	82.42 $\pm$ 0.6	9.3
<b>GT</b>	69.34 $\pm$ 1.2	55.04 $\pm$ 1.9	36.29 $\pm$ 1.0	<b>84.08 <math>\pm</math> 5.6</b>	80.00 $\pm$ 4.9	<b>84.80 <math>\pm</math> 4.3</b>	84.38 $\pm$ 1.3	91.28 $\pm$ 1.1	<b>79.93 <math>\pm</math> 0.5</b>	79.04 $\pm$ 0.5	70.16 $\pm$ 0.8	82.09 $\pm$ 0.7	5.6
<b>FAGCN</b>	60.98 $\pm$ 2.3	42.20 $\pm$ 1.8	35.67 $\pm$ 0.9	77.00 $\pm$ 7.7	78.32 $\pm$ 6.3	82.41 $\pm$ 3.8	82.79 $\pm$ 2.7	91.99 $\pm$ 1.0	79.27 $\pm$ 0.6	79.19 $\pm$ 0.4	71.55 $\pm$ 0.8	83.88 $\pm$ 0.5	7.5
<b>DMP</b>	63.79 $\pm$ 4.1	34.19 $\pm$ 7.6	28.30 $\pm$ 2.7	66.08 $\pm$ 7.0	56.41 $\pm$ 5.5	62.73 $\pm$ 4.5	70.58 $\pm$ 11.3	82.63 $\pm$ 4.1	56.41 $\pm$ 7.8	70.07 $\pm$ 4.1	59.12 $\pm$ 4.4	75.05 $\pm$ 3.8	12.8
<b>MixHop</b>	60.30 $\pm$ 2.1	41.05 $\pm$ 2.0	<b>36.48 <math>\pm</math> 1.2</b>	77.73 $\pm$ 7.3	75.95 $\pm$ 5.7	82.12 $\pm$ 4.5	79.52 $\pm$ 2.1	89.45 $\pm$ 1.5	78.59 $\pm$ 0.7	80.10 $\pm$ 0.4	71.42 $\pm$ 0.9	81.61 $\pm$ 0.8	9.3
<b>GPRGNN</b>	66.92 $\pm$ 1.7	46.32 $\pm$ 1.5	35.58 $\pm$ 0.9	81.51 $\pm$ 6.6	76.86 $\pm$ 7.1	84.06 $\pm$ 5.2	85.82 $\pm$ 0.9	<b>92.41 <math>\pm</math> 0.7</b>	79.67 $\pm$ 0.5	80.28 $\pm$ 0.4	71.59 $\pm$ 0.8	84.20 $\pm$ 0.5	<b>5.2</b>
<b>DAGNN</b>	54.99 $\pm$ 2.0	40.03 $\pm$ 1.4	33.69 $\pm$ 1.0	61.35 $\pm$ 6.1	63.89 $\pm$ 7.0	62.27 $\pm$ 4.2	<b>85.83 <math>\pm</math> 0.8</b>	92.30 $\pm$ 0.7	79.31 $\pm$ 0.6	<b>80.44 <math>\pm</math> 0.5</b>	73.16 $\pm$ 0.6	<b>84.43 <math>\pm</math> 0.5</b>	7.2
<b>AERO-GNN</b>	<b>71.58 <math>\pm</math> 2.4</b>	<b>61.76 <math>\pm</math> 2.4</b>	<b>36.57 <math>\pm</math> 1.1</b>	<b>84.35 <math>\pm</math> 5.2</b>	<b>81.24 <math>\pm</math> 6.8</b>	<b>84.80 <math>\pm</math> 3.3</b>	<b>86.69 <math>\pm</math> 1.4</b>	<b>92.50 <math>\pm</math> 0.7</b>	<b>79.95 <math>\pm</math> 0.5</b>	<b>80.59 <math>\pm</math> 0.5</b>	<b>73.20 <math>\pm</math> 0.6</b>	83.90 $\pm$ 0.5	<b>1.4</b>

- In each column, ■ indicates ranking the first, and ■ indicates ranking the second. A.R. denotes average ranking.



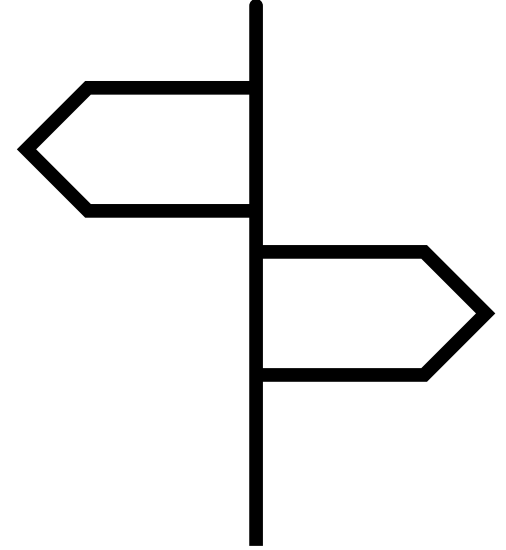
# Performance Over Layers

- **AERO-GNN** has
  - Highest best performance across model depth (see ★ in the Figure)
  - Better performance over layers  $k$  (see **trend** in the Figure)

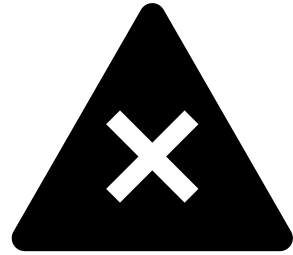


# Contents

- Sec. 1: Introduction
- Sec. 2: Analysis of Graph Attention
- Sec. 3: Proposed Method : AERO-GNN
- Sec. 4: Experiments and Empirical Evaluation
- Sec. 5: Discussion



# Summary



## **Problem**

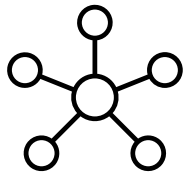
Two Limitations to  
Deep Graph Attention



## **Solution: AERO-GNN**

Theoretically and Empirically  
Mitigates the Problems

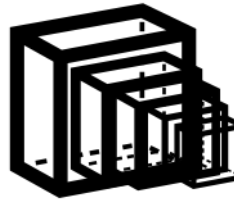
# Implications for Graph Learning



## Attention-Based GNNs

A larger focus has been placed on designing a **more expressive layer**

- *with new designs*
- *with new loss terms*
- *with more features*



## Deep GNNs

Making deeper GNNs have been an important **setback to GNN research**

- *over-smoothing*
- *over-squashing*
- *over-correlation*



## We Bridge the Two

The two are complementary

**Thank You**