



# Advanced Deep Graph Learning: Deeper, Faster, Robuster, and Unsupervised

Yu Rong, Tingyang Xu, Yatao Bian, Junzhou Huang, Tencent AI Lab  
Wenbing Huang, Fuchun Sun, Tsinghua University  
Hong Cheng, The Chinese University of Hong Kong

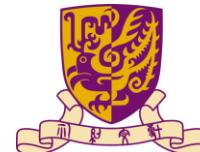
Tutorial website: <https://ai.tencent.com/ailab/ml/WWW-Deep-Graph-Learning.html>



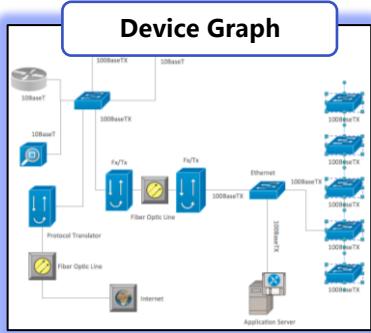
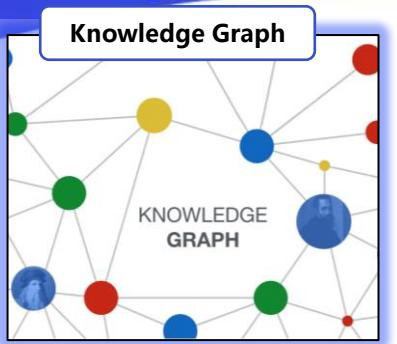
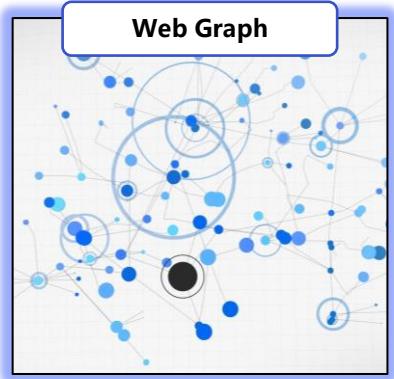
Tencent  
AI Lab



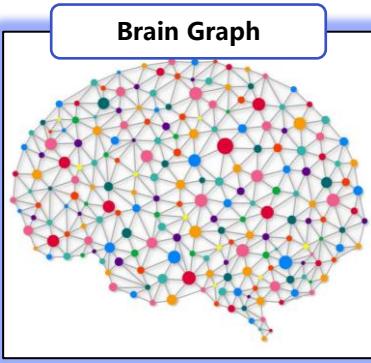
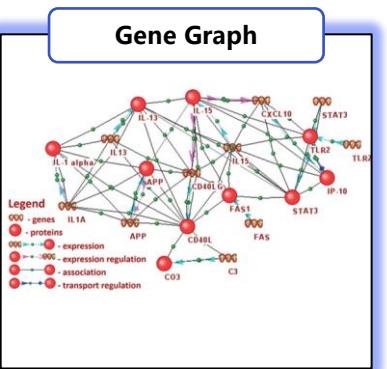
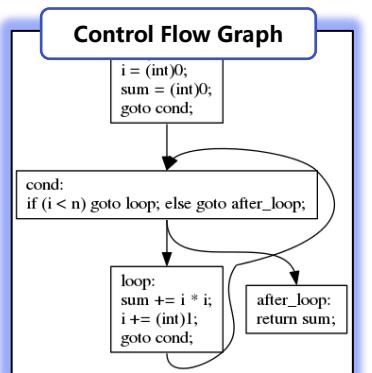
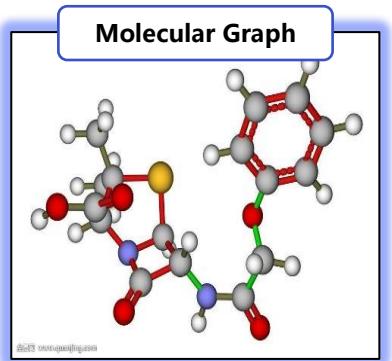
清华大学  
Tsinghua University



香港中文大學  
The Chinese University of Hong Kong

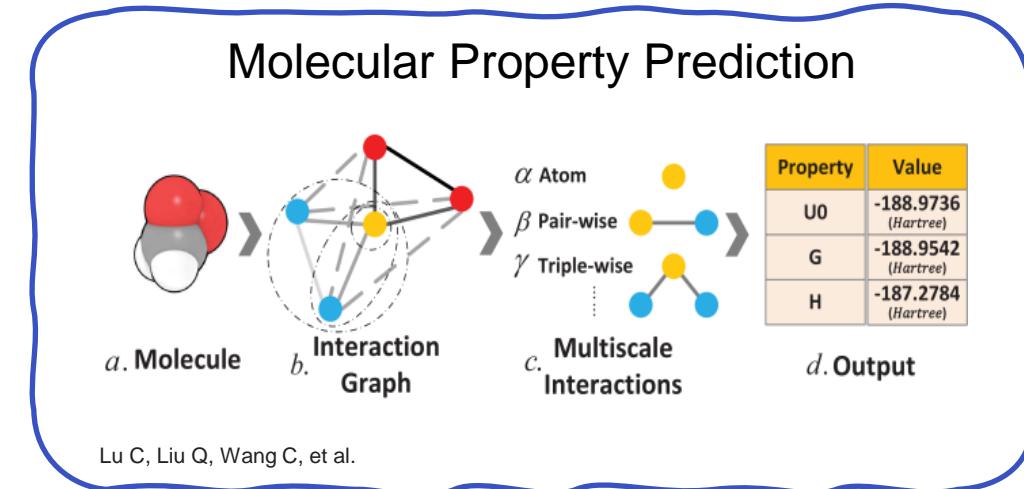
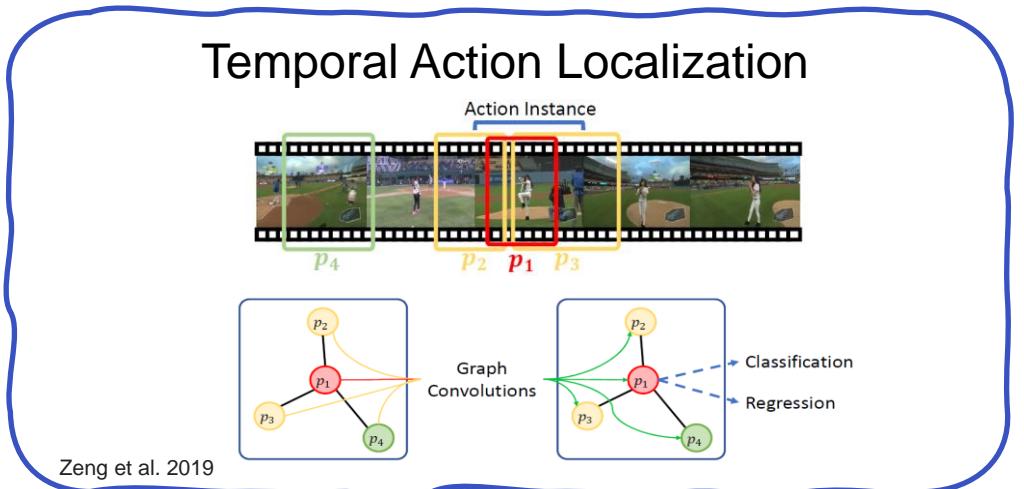
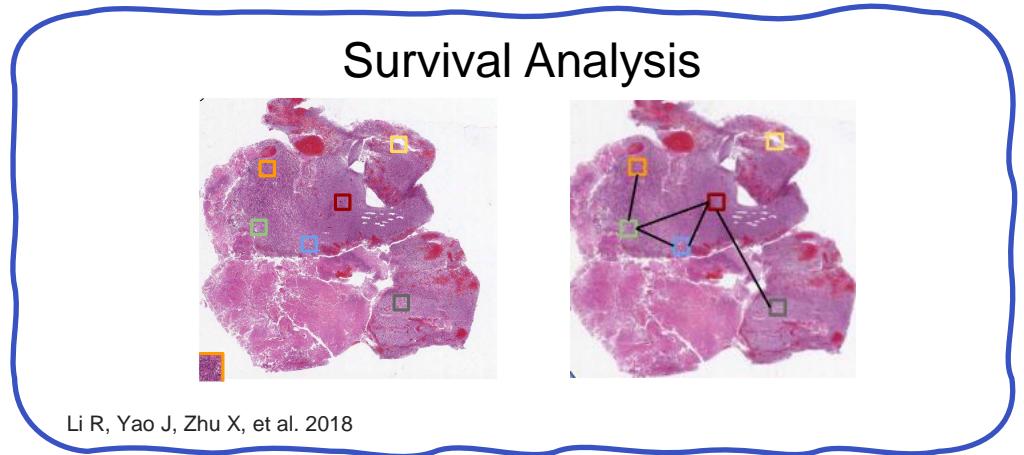
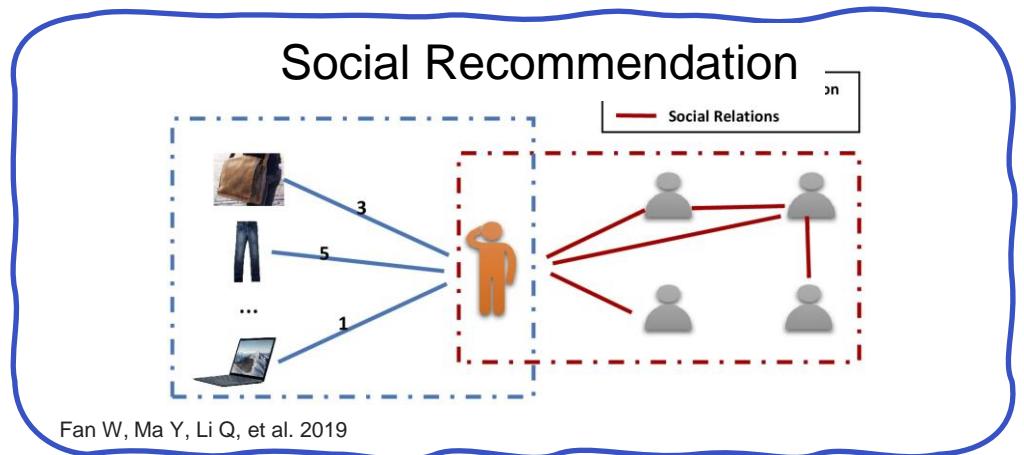


# Graph is Everywhere



# Graph is Important

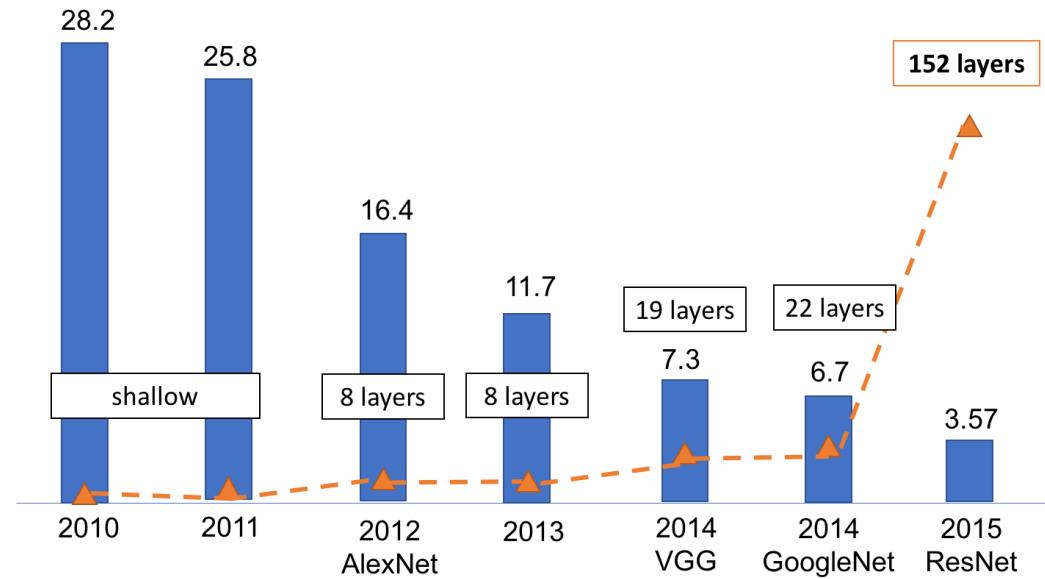
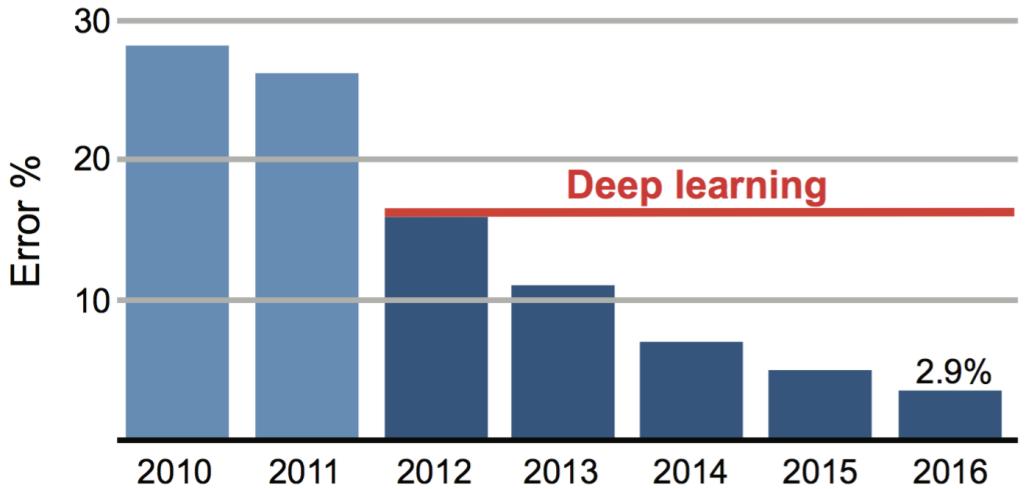
- Numerous real-world problems can be summarized as a set of tasks on graphs.





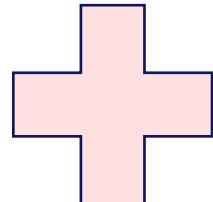
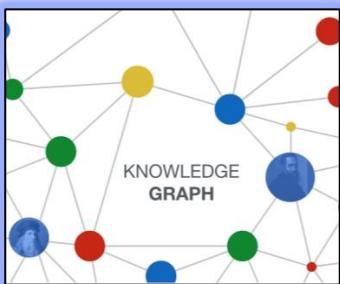
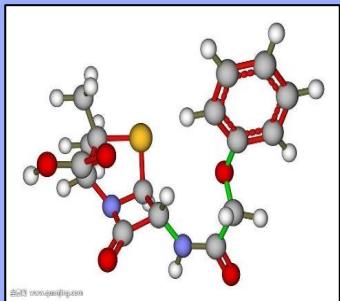
# The Power of Deep Learning

IMAGENET

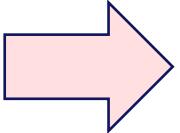
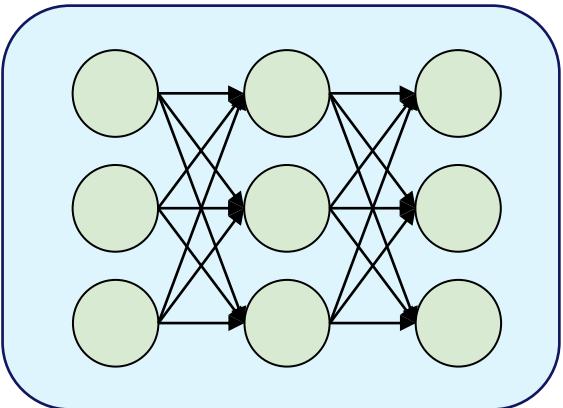




# Deep Learning + Graphs = ?

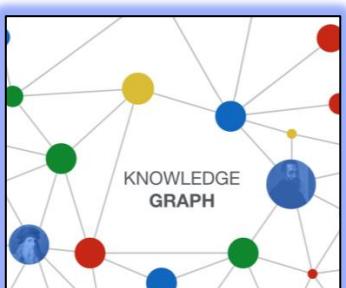
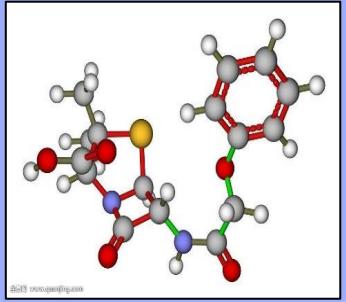


Deep Learning Model

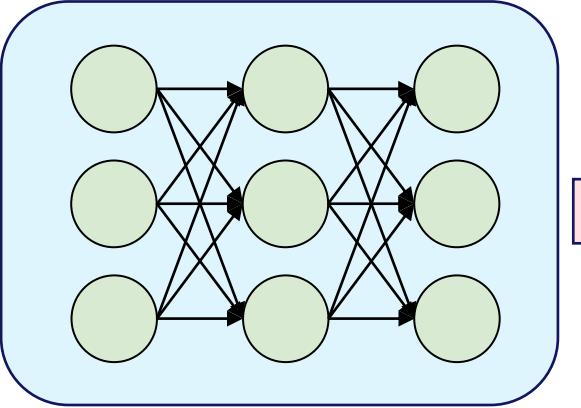


?

# Deep Graph Learning

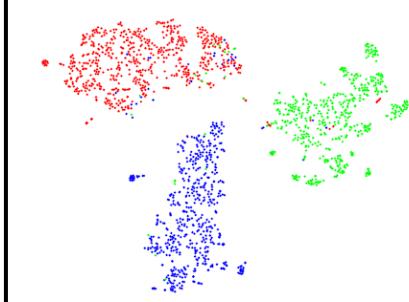


Neural Network



Make neural network model that can deal with graph data.

Graph/Node Representation



## Applications

Node Classification

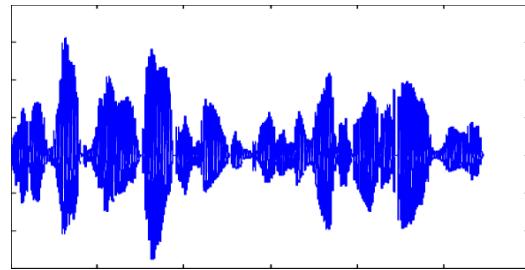
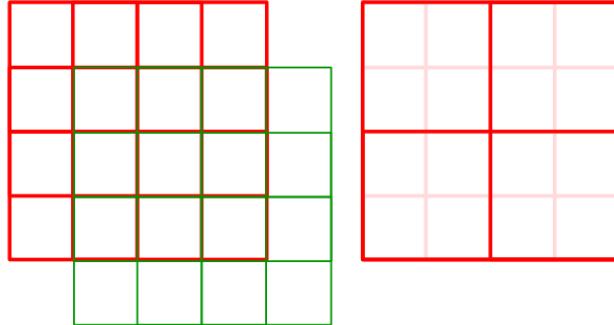
Link Prediction

Community Detection

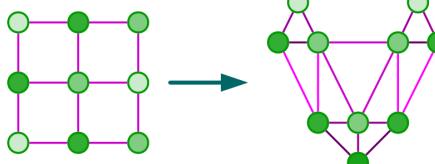
Graph Generation

.....

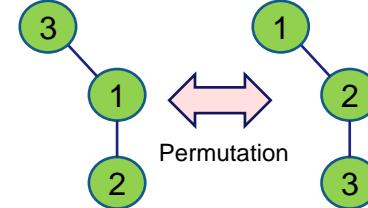
# The Big Challenges



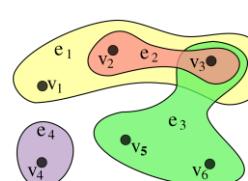
**VS**



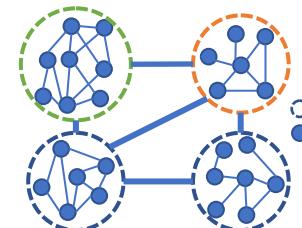
Irregular



Permutation  
Invariant

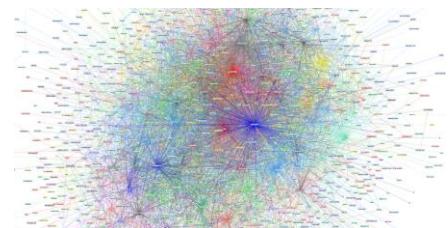


HyperGraph



Hierarchical Graph

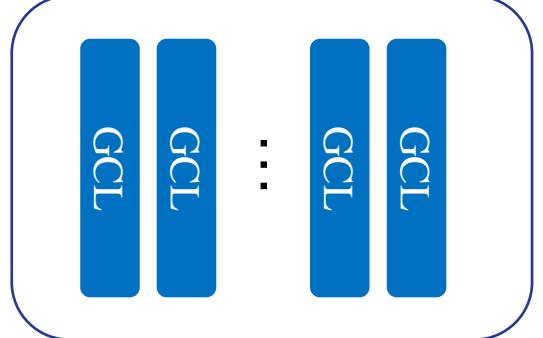
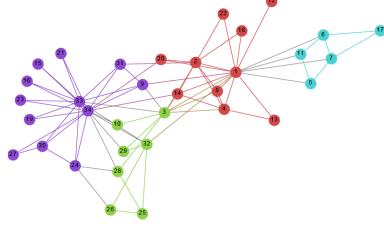
Complex variants



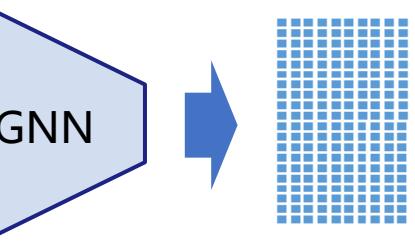
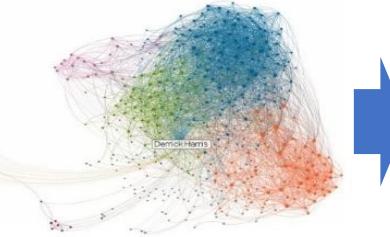
Large-scale instance

- Grid-like and sequence-like structure.
- Spacial/sequential relations between pixels / units.

# The Research Questions

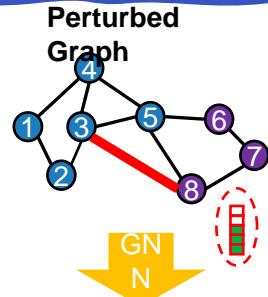
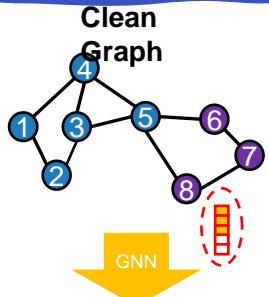


How to train the deeper graph neural networks?

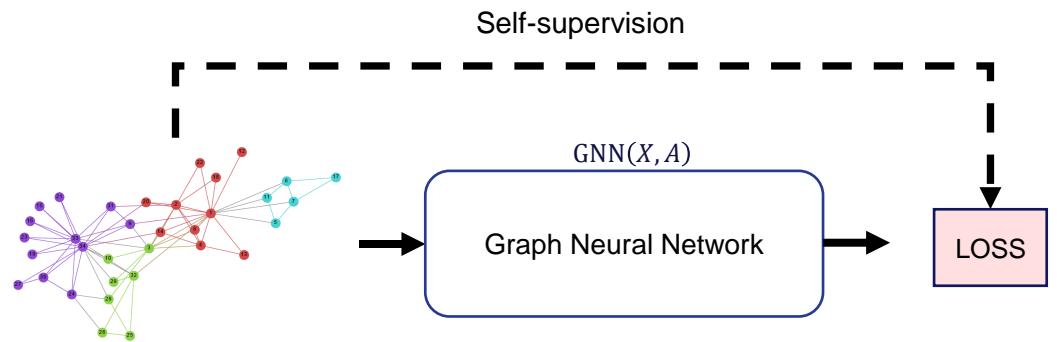


Node  
Representation

How to extend the graph neural networks to process a large-scale graph?



How to enhance the robustness of graph neural network?



How to conduct the self/un-supervised learning on graphs?

# Tutorial Overview

## Foundations

- Preliminary
- The Brief History of Graph Neural Networks

## Advances

- Training Deep GNNs
- Scalability of GNNs
- Robustness of GNNs
- Self/Un-Supervised Learning of GNNs
- Other Advanced Topics

## Applications

- Social Networks
- Medical Imaging

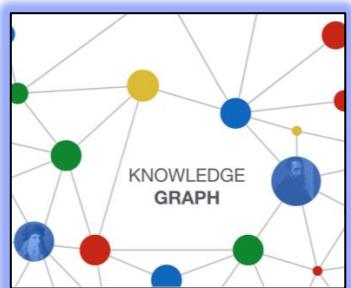


Tencent  
AI Lab

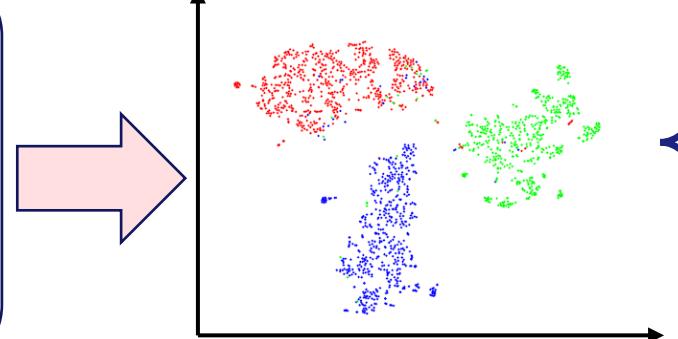
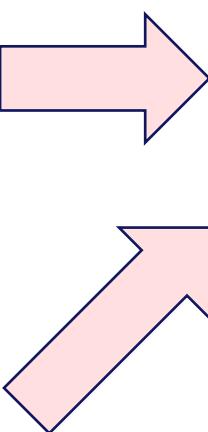


# **Preliminaries and Brief History of Graph Neural Networks**

# What is the Graph Neural Network?

Graph Neural Network



Graph/Node  
Representation

Applications

Node  
Classification

Link Prediction

Community  
Detection

Graph  
Generation

.....

Neural network model that can deal with graph data.

# Graph Neural Network is not a New Thing

Sperduti, Alessandro and Starita, Antonina. 1997

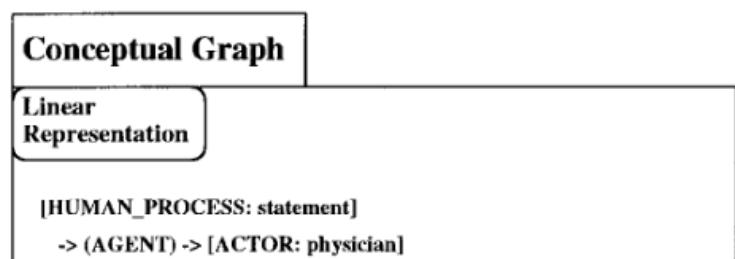
714

IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 8, NO. 3, MAY 1997

## Supervised Neural Networks for the Classification of Structures

Alessandro Sperduti and Antonina Starita, *Member, IEEE*

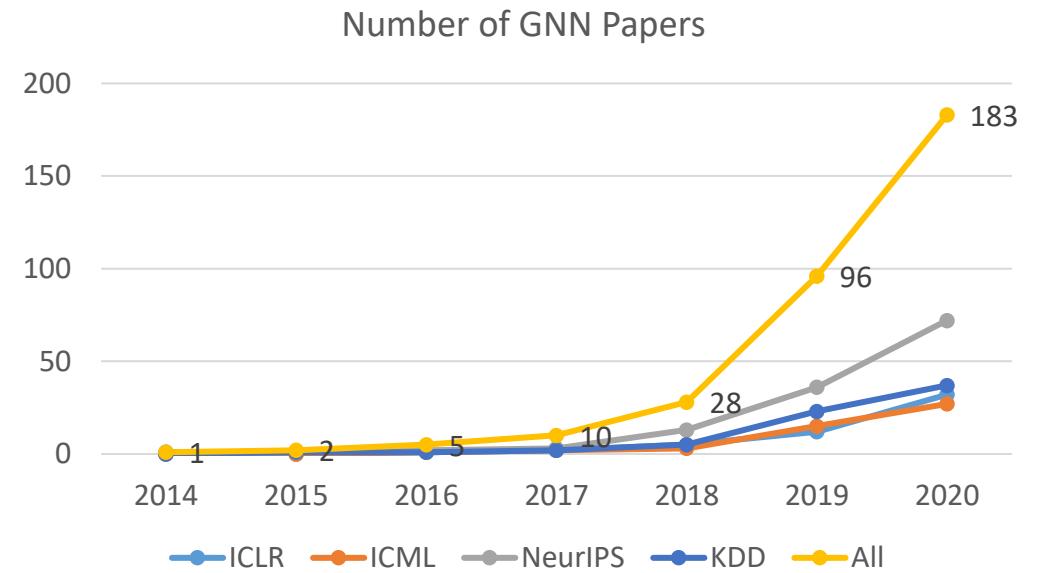
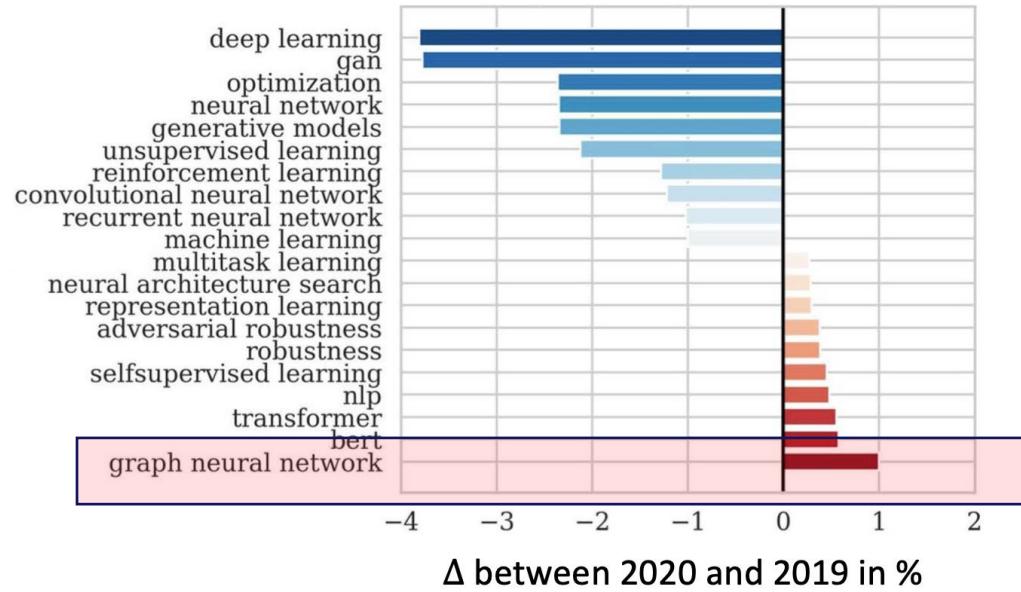
**Abstract**—Until now neural networks have been used for classifying unstructured patterns and sequences. However, standard neural networks and statistical methods are usually believed to be inadequate when dealing with complex structures because of their feature-based approach. In fact, feature-based approaches usually fail to give satisfactory solutions because of the sensitivity of the approach to the *a priori* selection of the features, and the incapacity to represent any specific information on the relationships among the components of the structures.





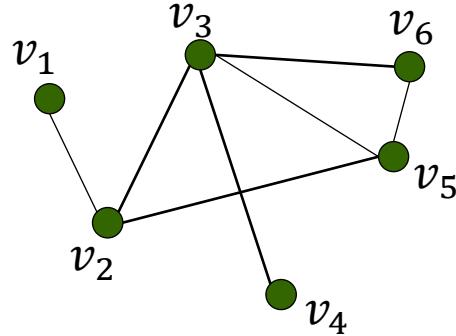
# A Rapidly Growing Area

ICLR 2020 submissions keyword statistics



# Preliminaries of Graph Learning

A topological graph



$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$$

which has n nodes  
and m edges

$$\mathcal{V} = \{v_1, \dots, v_n\}$$

$$\mathcal{E} = \{e_1, \dots, e_m\}$$

Node features:  
 $X \in R^{n \times d}$

**Adjacent matrix  $A$** :  $A_{ij} = 1$ , existing edge between  $i$  and  $j$   
 $A_{ij} = 0$ , not exist edge between  $i$  and  $j$

**Degree matrix  $D$**  =  $\text{diag}(\text{degree}(v_1), \dots, \text{degree}(v_n))$

**Laplacian matrix  $L$**  =  $D - A$

Degree matrix

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{pmatrix}$$

Adjacency matrix

$$- \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{pmatrix} =$$

Laplacian matrix

$$\begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 4 & -1 & -1 & -1 \\ 0 & 0 & -1 & 1 & 0 & 0 \\ 0 & -1 & -1 & 0 & 3 & -1 \\ 0 & 0 & -1 & 0 & -1 & 2 \end{pmatrix}$$

# The Model of Graph Neural Networks



# The Model of Graph Neural Networks

GNN 1.0

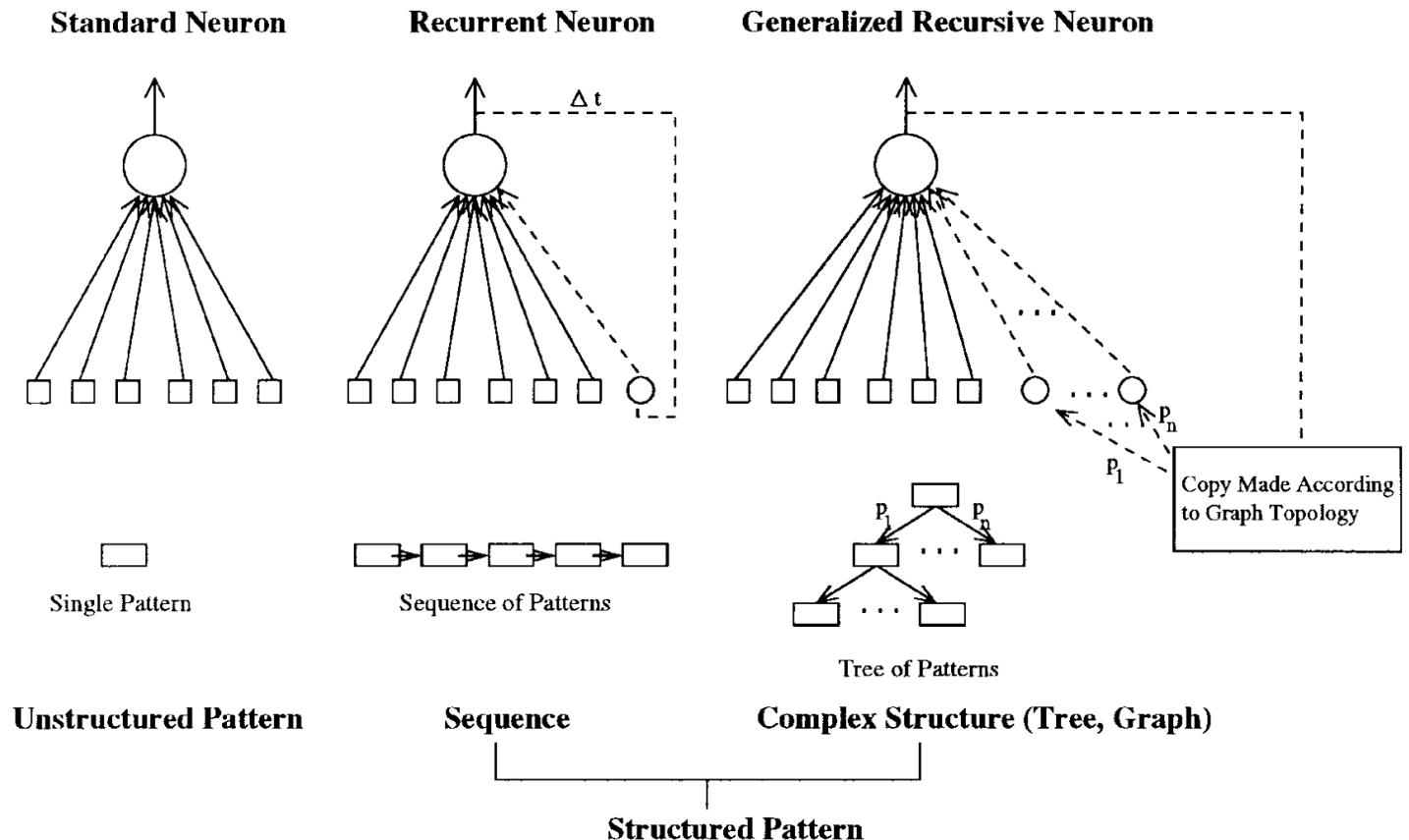
- Understanding GNN as RNN

GNN 2.0

GNN 3.0

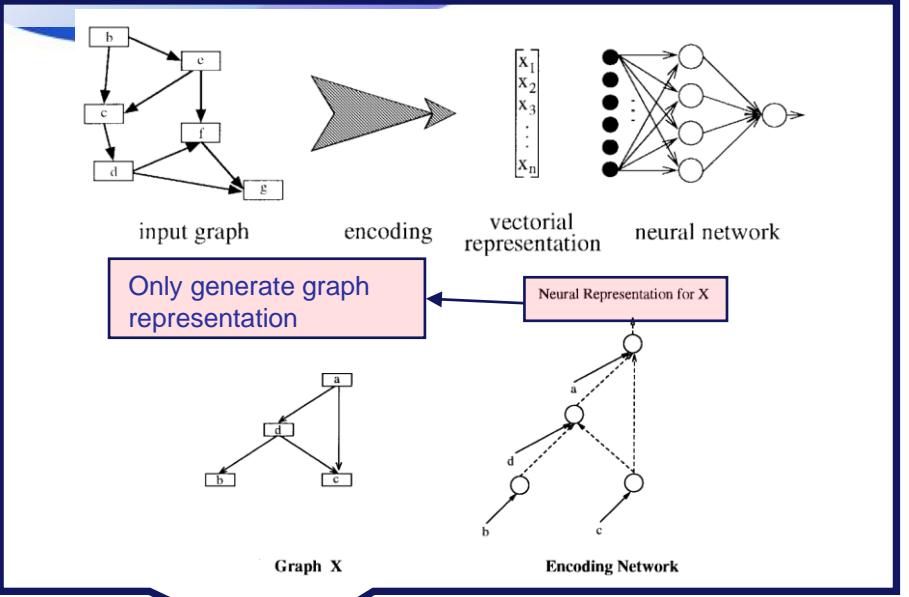


# GNN 1.0: Understanding GNN as RNN



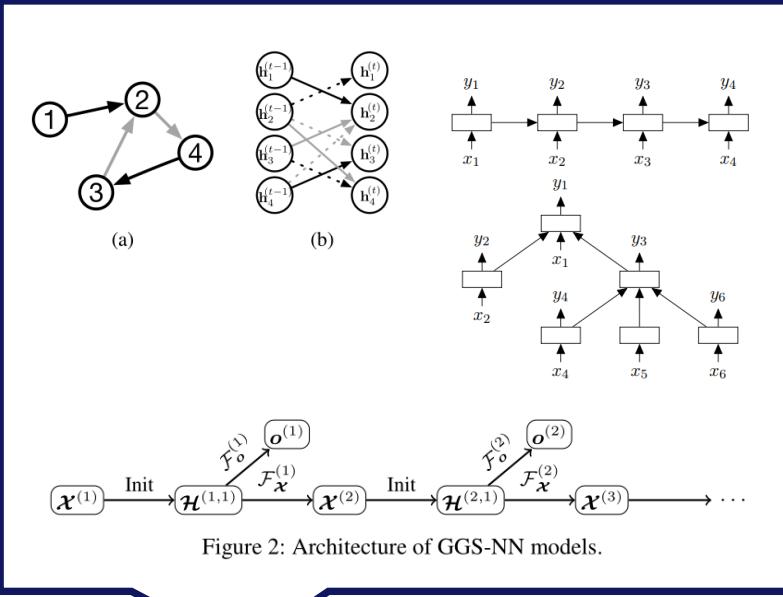
- The RNN on sequences can be generalized to trees and DAGs.

# GNN 1.0: Understanding GNN as RNN



## From 2000 to 2010

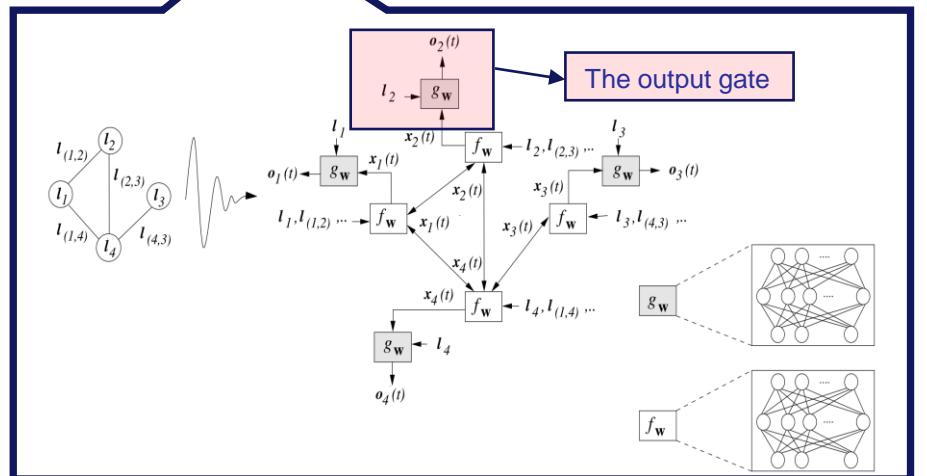
Gori et.al (IJCNN 05) and Scarselli et.al (TNN 08) add **the output gate** for each node to generate the node representation in graphs. This model is called GraphRNN.



## Before 2000

Sperduti, Alessandro, and Antonina Starita. (TNN 97) propose the **generalized recursive neuron** for the graph classification problem on Trees/DAGs.

This generalized recursive neuron can only generate the graph representations.



## After 2010

Li, Yujia, et al. (ICLR 16) add gated recurrent units and modern optimization techniques to improve the performance of Scarselli et.al (TNN 09).

Tai, Kai Sheng et.al. (ACL 2015) extend LSTM to a tree-structured network topologies.

# The Brief History of Graph Neural Networks

GNN 1.0

- Understanding GNN as RNN

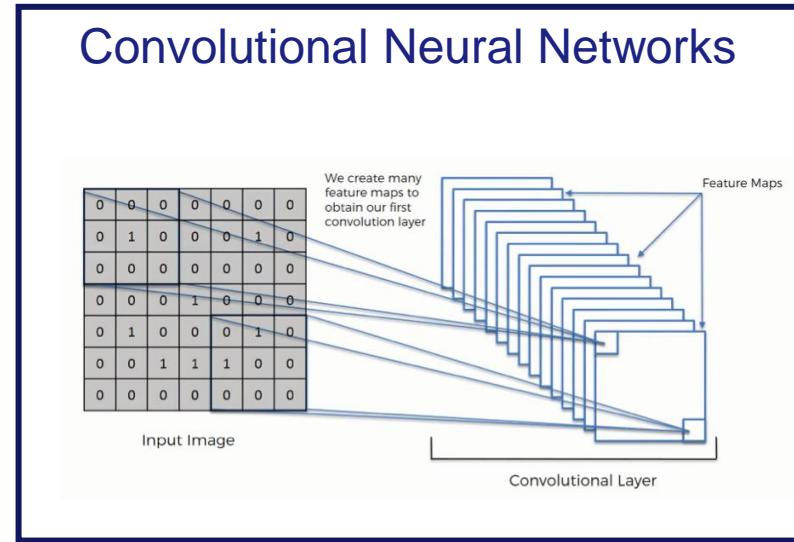
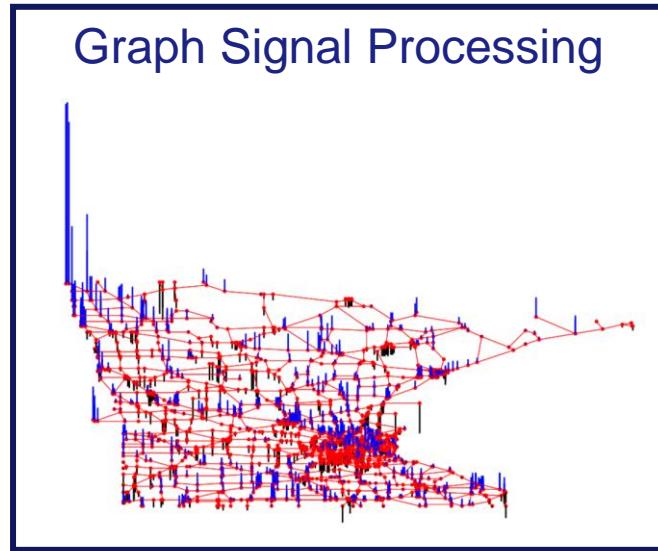
GNN 2.0

- Understanding GNN as Convolution

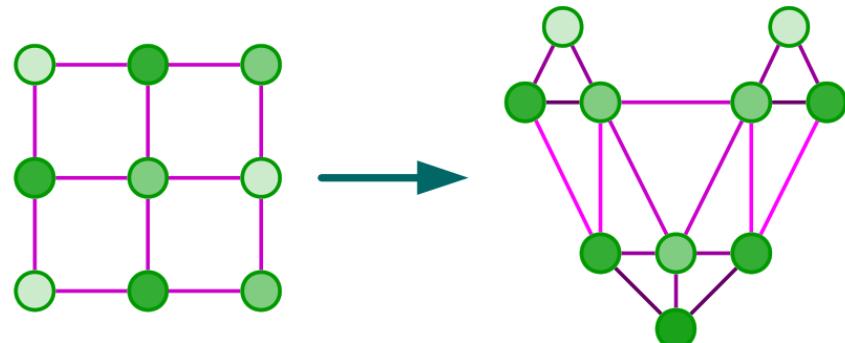
GNN 3.0



# GNN 2.0: Understanding GNN as Convolution

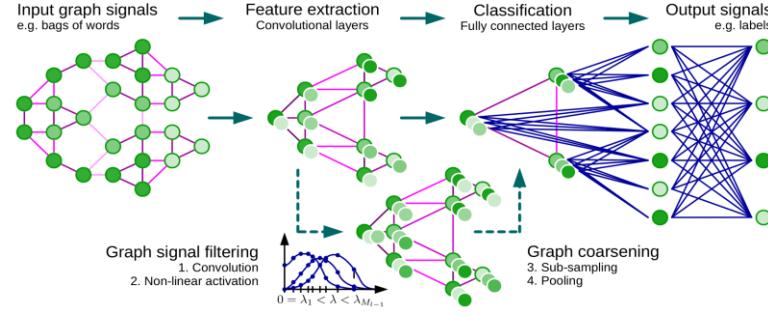


- How to perform the convolutions on graphs?
  - Irregular structures.
  - Weighted edges.
  - No orientation or ordering (in general).



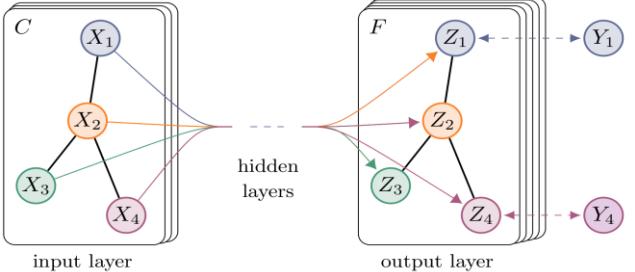
# GNN 2.0: Understanding GNN as Convolution

## ChebNet (NeurIPS 2016) [2]



- Build the connection between graph signal processing and graph convolution.
- Use Chebyshev polynomial to fast approximate the graph filtering in the spectral domain.

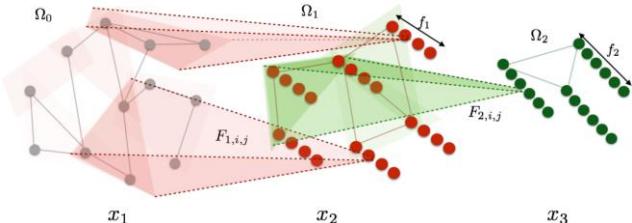
## Graph Convolutional Network (ICLR 2017)



$$H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}H^{(l)}W^{(l)})$$

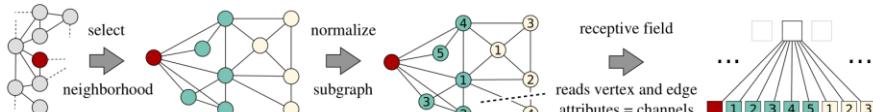
- Approximate 1-order Chebyshev polynomial the in spatial domain.
- Layer-wise convolution to extend receptive field.
- **The practical convolutional model for graphs.**

## Deep Locally Connected Networks (ICLR 2014) [1]



- Discuss two constructions on both spatial and spectral domain.
- Analog the convolution operation based on the Laplacian spectrum.
- **Additional eigen decomposition is needed.**

## PATCHY-SAN (ICML 2016)

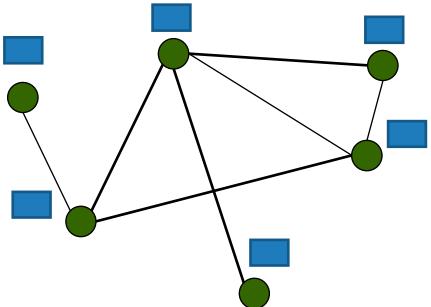


- Neighborhood sampling to construct receptive field.

# Graph Signal Processing

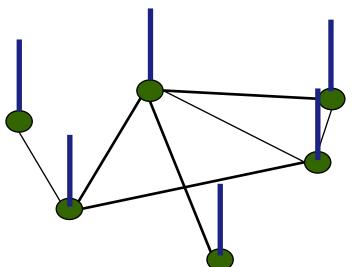
Graph signal:

$$h: \mathcal{V} \rightarrow \mathbb{R}^n$$

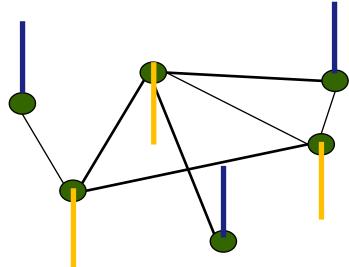


**“Frequency” or “Smoothness” of the signal  $h$**

$$h^T L h = \sum_{i < j} A_{ij} (h_i - h_j)^2$$



Low frequency graph signal



High frequency graph signal

Eigen decomposition of graph Laplacian

$$L = U \Lambda U^T$$

$$L = \begin{bmatrix} | & & | \\ \mathbf{u}_0 & \cdots & \mathbf{u}_{N-1} \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_0 & & 0 \\ & \ddots & \\ 0 & & \lambda_{N-1} \end{bmatrix} \begin{bmatrix} \mathbf{u}_0 & & \\ \vdots & \ddots & \\ \mathbf{u}_{N-1} & & \end{bmatrix}$$

eigenvalues sorted non-decreasingly:

$$0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_{n-1}$$

The frequency of an eigenvector of  $L$  is its corresponding eigenvalue:

$$u_i^T L u_i = u_i^T \lambda_i u_i = \lambda_i$$



# Graph Convolution: Spectral domain → Spatial domain

**Graph Convolution:** input signal  $x$ , filter  $\hat{g}$ , graph Laplacian  $L$

Spectral

$$y = x *_{\mathcal{G}} \hat{g} = \mathbf{U} \begin{pmatrix} \hat{g}(\lambda_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{g}(\lambda_n) \end{pmatrix} \mathbf{U}^T x = \mathbf{U} \hat{g}(\Lambda) \mathbf{U}^T x =: \hat{g}(L) x$$

**Parameterization:** replace  $\hat{g}(\Lambda)$  with  $\hat{g}_\theta = \text{diag}(\theta)$

**ChebNet** (NeurIPS 2016): parameterize with Chebyshev polynomials:

$$y = \hat{g}_\theta(L)x = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{L})x \quad \tilde{L} = \frac{1}{\lambda_{max}} L - I$$

**GCN** (ICLR 2017): simplified ChebNet  $K = 1$ , suppose  $\lambda_{max} = 2$ ,  $\theta := \theta_0 = -\theta_1$

$$\hat{g}_\theta(L)x = (I + D^{-\frac{1}{2}} A D^{-\frac{1}{2}}) x \theta$$



Spatial

Apply a renormalization trick:  $\hat{g}_\theta(L)x = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} x \theta \quad \tilde{A} = A + I$  (add self-loop)

# The Brief History of Graph Neural Networks

GNN 1.0

- Understanding GNN as RNN

GNN 2.0

- Understanding GNN as Convolution

GNN 3.0

- Variants of Convolutions
- GNN with Attention
- GNN with Graph Pooling



# GNN 3.0: Variants of Convolutions

$$g_\theta * x = U g_\theta U^T x \quad \longrightarrow \quad \mathbf{H}^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)})$$

Lanczos Network [3]

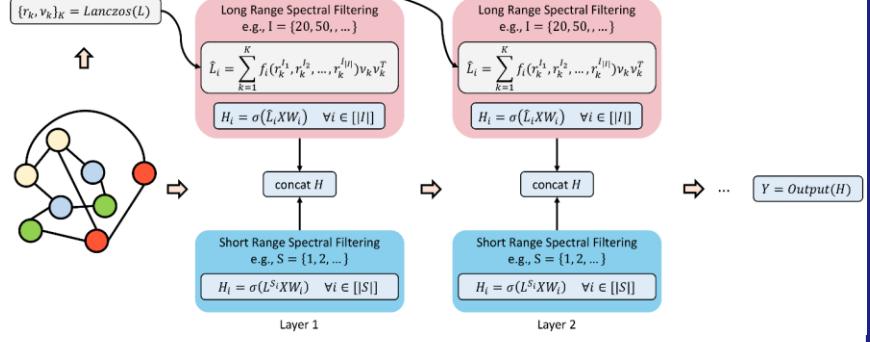
Graph Wavelet Neural Network [1]

Hyperbolic GCN [2]

# GNN 3.0: Variants of Convolutions

$$g_\theta * x = U g_\theta U^T x \quad \longrightarrow \quad H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

Lanczos Network [3]



Graph Wavelet Neural Network [1]



Hyperbolic GCN [2]

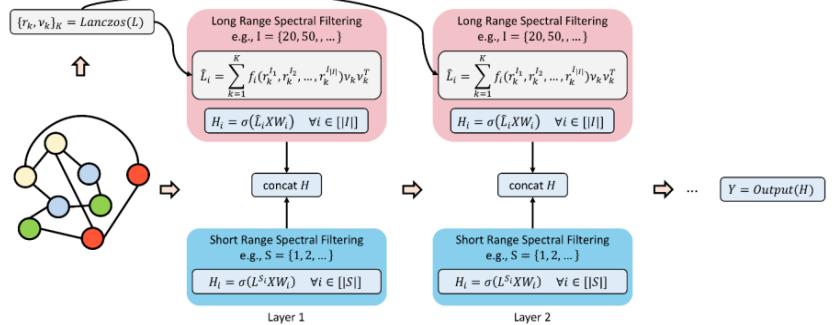


- Employ Lanczos algorithm to obtain the low-rank approximation of the graph Laplacian  $I - \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ .
- Easy to construct multi-scale Graph Convolution.

# GNN 3.0: Variants of Convolutions

$$g_\theta * x = \boxed{U g_\theta U^T x} \quad \longrightarrow \quad H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

## Lanczos Network [3]



- Employ Lanczos algorithm to obtain the low-rank approximation of the graph Laplacian  $I - \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ .
- Easy to construct multi-scale Graph Convolution.

## Graph Wavelet Neural Network [1]

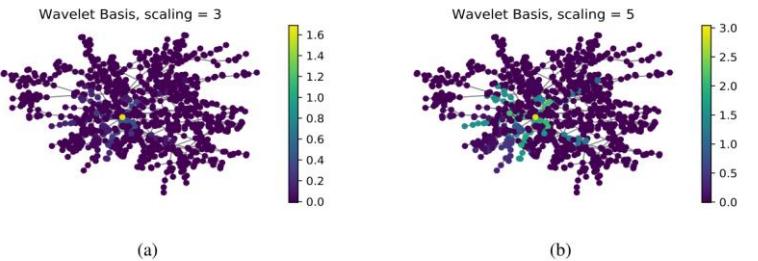


Figure 1: Wavelets on an example graph at (a) small scale and (b) large scale.

$$H_{[:,j]}^{(l+1)} = \sigma \left( \psi_s \sum_{i=1}^p F_{i,j}^{(l)} \psi_s^{-1} H_{[:,i]}^{(l)} \right), \quad j = 1, \dots, q$$

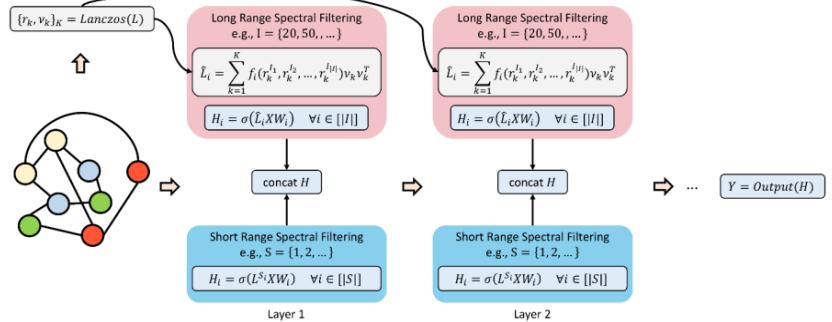
- Use wavelet transform to replace Fourier transform in the original GCN.
- More localized convolution and flexible neighborhood.

## Hyperbolic GCN [2]

# GNN 3.0: Variants of Convolutions

$$g_\theta * x = U g_\theta U^T x \quad \longrightarrow \quad H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

## Lanczos Network [3]



- Employ Lanczos algorithm to obtain the low-rank approximation of the graph Laplacian  $I - \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ .
- Easy to construct multi-scale Graph Convolution.

## Graph Wavelet Neural Network [1]

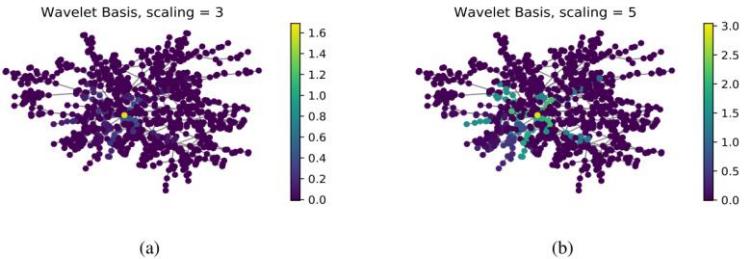
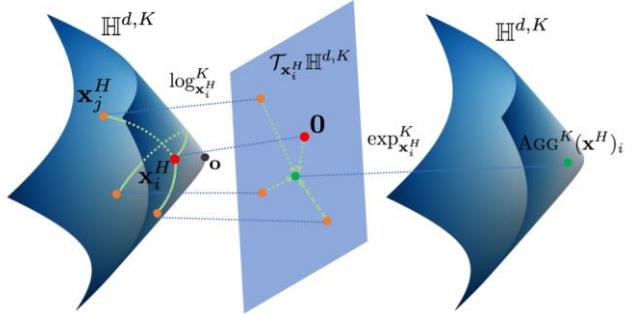


Figure 1: Wavelets on an example graph at (a) small scale and (b) large scale.

$$H_{[:,j]}^{(l+1)} = \sigma \left( \psi_s \sum_{i=1}^p F_{i,j}^{(l)} \psi_s^{-1} H_{[:,i]}^{(l)} \right), \quad j = 1, \dots, q$$

- Use wavelet transform to replace Fourier transform in the original GCN.
- More localized convolution and flexible neighborhood.

## Hyperbolic GCN [2]



Construct the GCN in hyperbolic space.

- Smaller distortion.
- Suitable for scale-free and hierarchical structure.
- Hyperbolic feature transform.  $h_i^{(l+1),H} = (W^{(l+1)} \otimes^{K_l} h_i^{(l),H}) \oplus^{K_l} b^{(l+1)}$
- Attention-based hyperbolic aggregation.  $y_i^{(l+1),H} = \text{AGG}^{K_l}(h^{(l),H})_i$

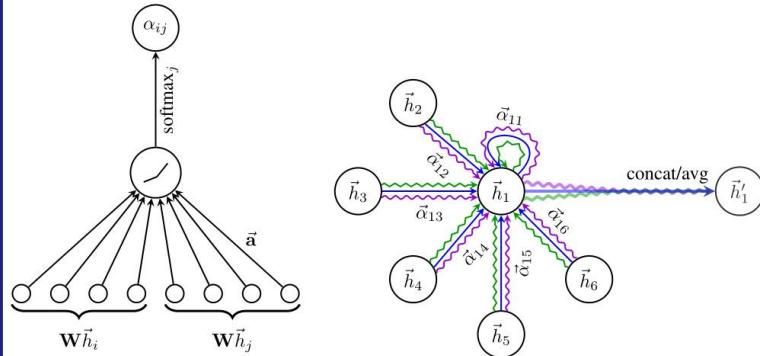
# GNN 3.0: GNN with Attention

$$\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)})$$

$$S = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$$

Fixed during training

Graph Attention Network [1]

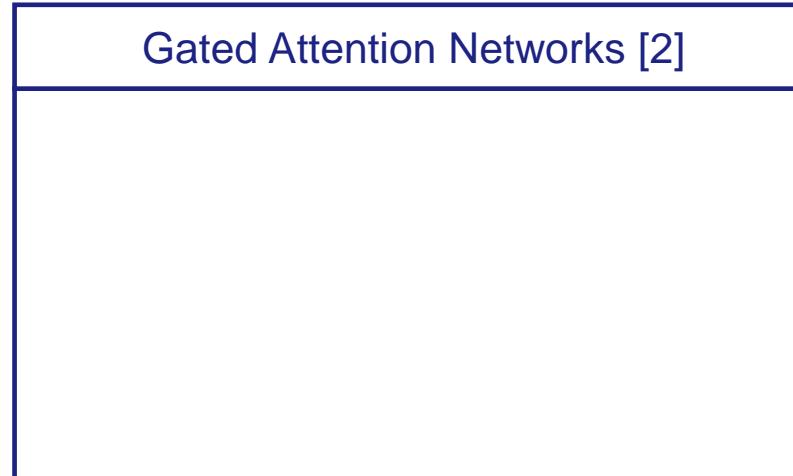


Replace the fixed aggregation weight  $a_{ij}$  to the learnable self-attention.

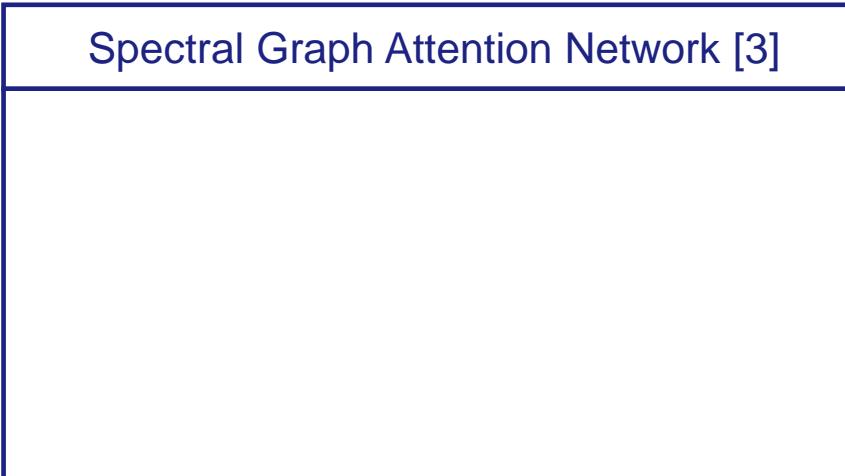
$$\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} a_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)})$$

$$a_{ij} = \exp\left(\frac{\sigma(\alpha^T [\mathbf{W}\mathbf{h}_i] || [\mathbf{W}\mathbf{h}_j])}{\sum_{k \in N(v_i)} \sigma(\alpha^T [\mathbf{W}\mathbf{h}_i] || [\mathbf{W}\mathbf{h}_k])}\right)$$

Gated Attention Networks [2]



Spectral Graph Attention Network [3]



# GNN 3.0: GNN with Attention

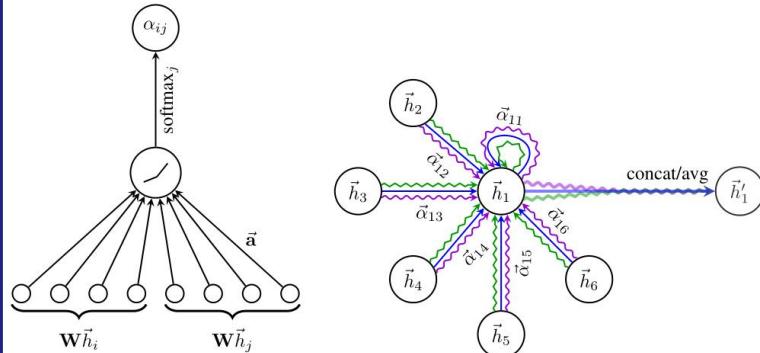
The original form:

$$\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)})$$

$$\mathbf{S} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$$

Fixed during training

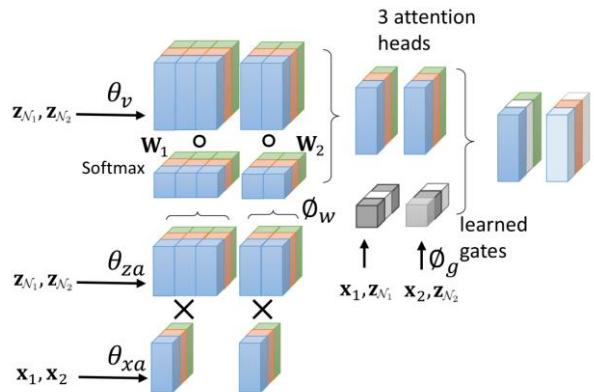
## Graph Attention Network [1]



Replace the fixed aggregation weight  $a_{ij}$  to the learnable self-attention.

$$\begin{aligned} \mathbf{h}_i^{(l+1)} &= \sigma\left(\sum_{j \in N(v_i)} \mathbf{a}_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}\right) \\ \mathbf{a}_{ij} &= \exp\left(\frac{\sigma(\alpha^T [\mathbf{W}\mathbf{h}_i] || [\mathbf{W}\mathbf{h}_j])}{\sum_{k \in N(v_i)} \sigma(\alpha^T [\mathbf{W}\mathbf{h}_i] || [\mathbf{W}\mathbf{h}_k])}\right) \end{aligned}$$

## Gated Attention Networks [2]



Add a learnable gate  $g_i^k$  to model the importance for each head.

$$\mathbf{h}_i^{(l+1)} = \sigma\left(\sum_{k=1}^K g_i^k \sum_{j \in N(v_i)} a_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}\right)$$

K is the number of heads.

## Spectral Graph Attention Network [3]

# GNN 3.0: GNN with Attention

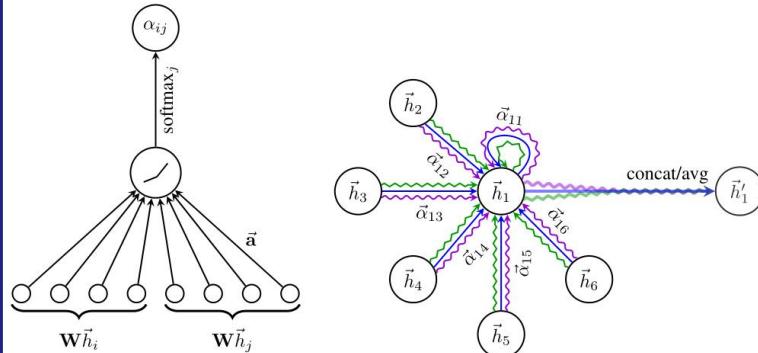
The original form:

$$\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)})$$

$$\mathbf{S} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$$

Fixed during training

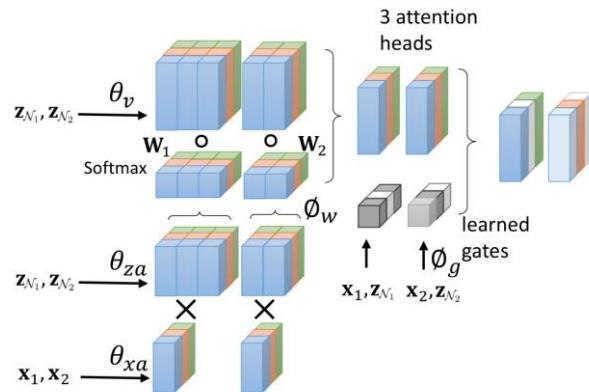
## Graph Attention Network [1]



Replace the fixed aggregation weight  $a_{ij}$  to the learnable self-attention.

$$\begin{aligned} \mathbf{h}_i^{(l+1)} &= \sigma\left(\sum_{j \in N(v_i)} \mathbf{a}_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}\right) \\ \mathbf{a}_{ij} &= \exp\left(\frac{\sigma(\alpha^T [\mathbf{Wh}_i] || [\mathbf{Wh}_j])}{\sum_{k \in N(v_i)} \sigma(\alpha^T [\mathbf{Wh}_i] || [\mathbf{Wh}_k])}\right) \end{aligned}$$

## Gated Attention Networks [2]

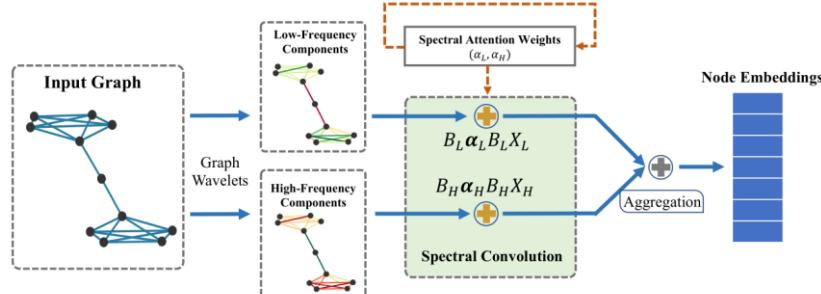


Add a learnable gate  $g_i^k$  to model the importance for each head.

$$\mathbf{h}_i^{(l+1)} = \sigma\left(\sum_{k=1}^K g_i^k \sum_{j \in N(v_i)} a_{ij} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}\right)$$

K is the number of heads.

## Spectral Graph Attention Network [3]



Apply the attention on the high / low-frequency components in spectral domain.

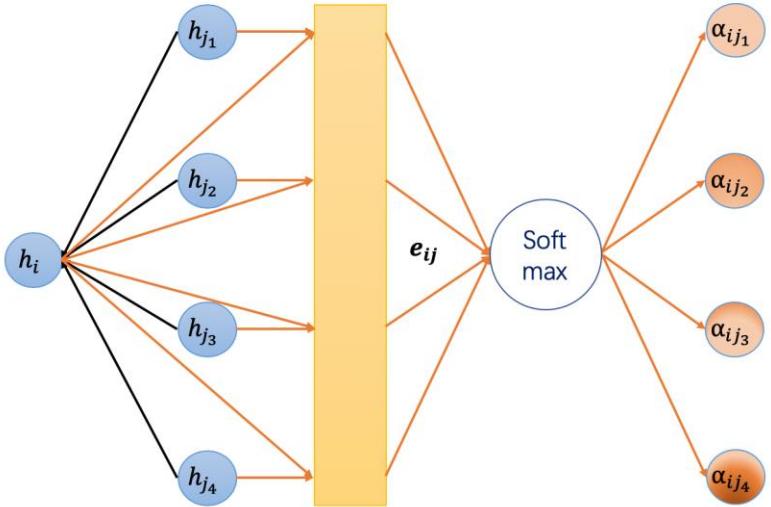
$$\mathbf{H}^{(l+1)} = \sigma(\text{AGG}(\mathbf{B}_L \mathbf{a}_L \mathbf{B}_L \mathbf{H}^{(l)}, \mathbf{B}_H \mathbf{a}_H \mathbf{B}_H \mathbf{H}^{(l)}) \mathbf{W}^{(l)})$$

$\mathbf{B} = [\mathbf{B}_L, \mathbf{B}_H]$  is the spectral graph wavelet bases.



# Graph Attention Network in Detail

## Single head attention



## Multi-head attention

Enrich the model capacity and stabilize the learning process

Each head has its own parameters and their outputs can be merged in two ways:

- Concatenation
- Average

$$\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} \alpha_{ij} \mathbf{W} \mathbf{h}_j^{(l)})$$

$$\alpha_{ij} = \text{Softmax}(e_{ij})$$

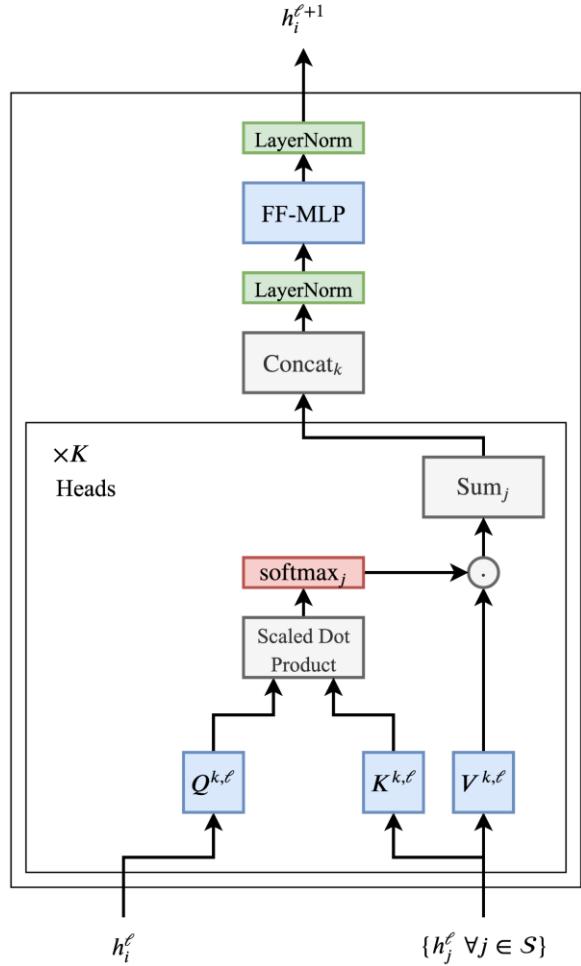
$$e_{ij} = \text{LeakyReLU}(\mathbf{a}^T (\mathbf{W}\mathbf{h}_i || \mathbf{W}\mathbf{h}_j))$$



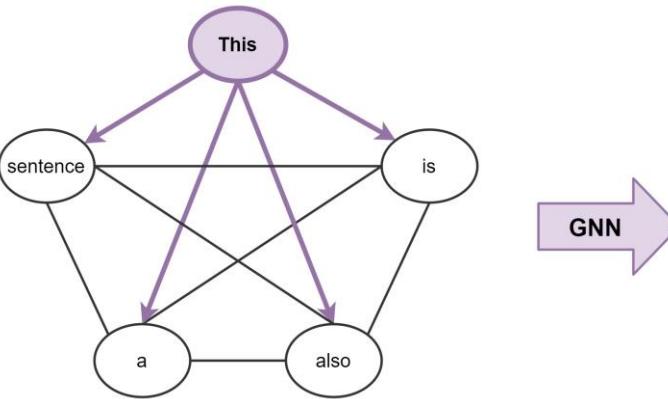
Attention weights learnt for the Cora dataset

# Transformers as GNNs with Multi-head Attention

One layer  
of the multi-  
head  
QKV attention



Transformer takes input sequence  
as a complete graph.



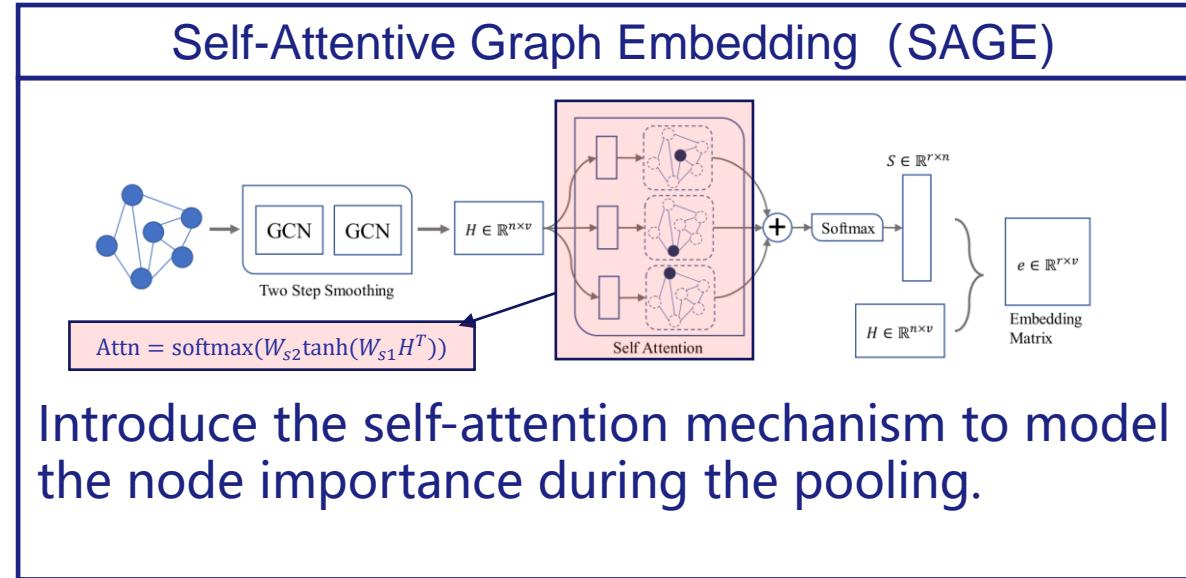
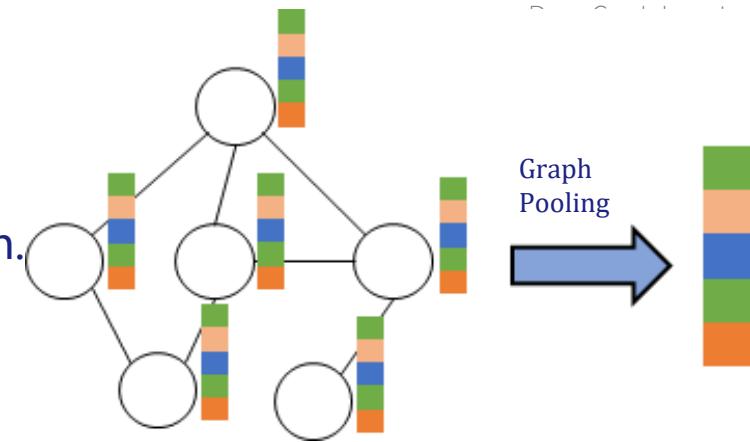
- Translation?
- Sentiment?
- Next word?
- Part-of-speech tags?

- ❑ Transformers can be viewed as GNNs with multi-head attention as the neighborhood aggregation function
- ❑ Transformers for NLP tasks treat the entire sequence as the neighborhood

# GNN 3.0: GNN with Graph Pooling

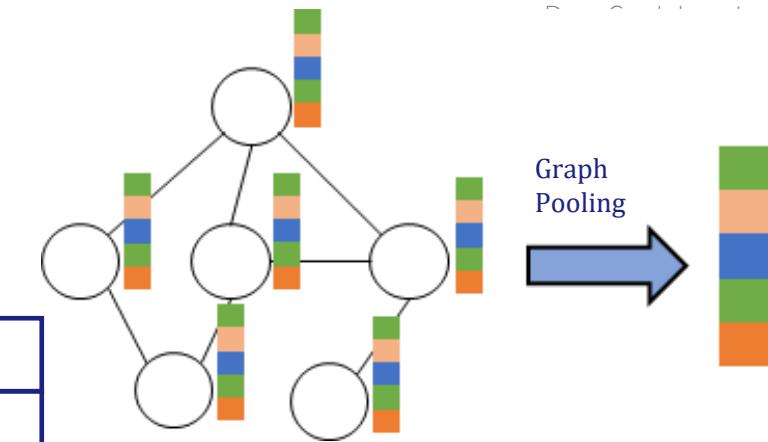
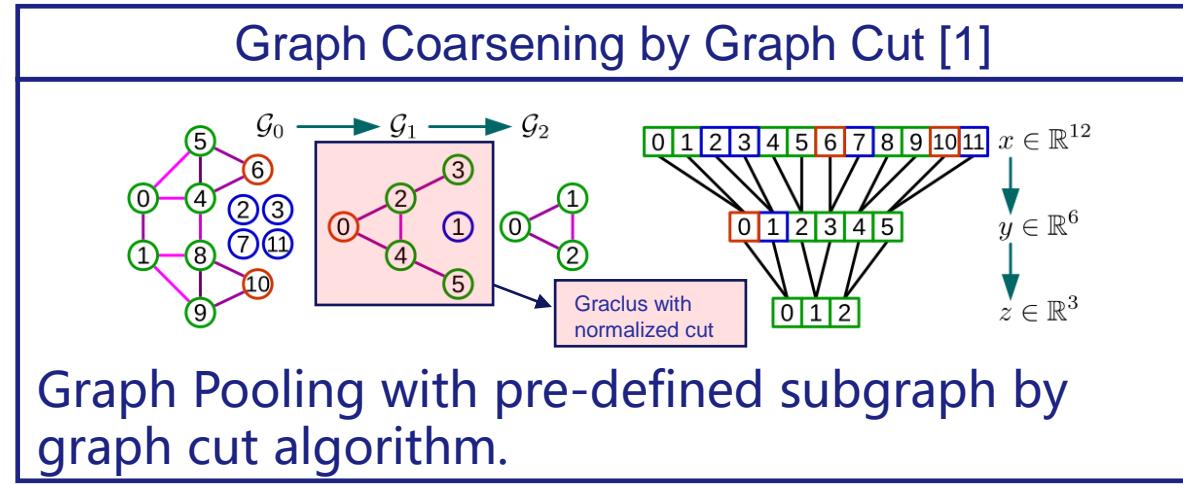
**Graph Pooling/Coarsening:** Convert the node representation to graph representation.

- The most straightforward way: Max/Mean Pooling
- SAGE: Attentive Pooling

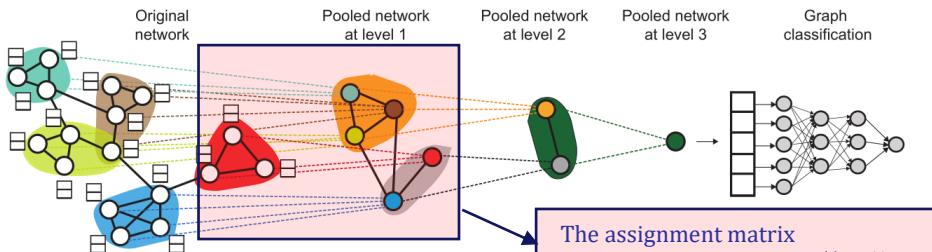


# GNN 3.0: GNN with Graph Pooling

## Hierarchical Pooling

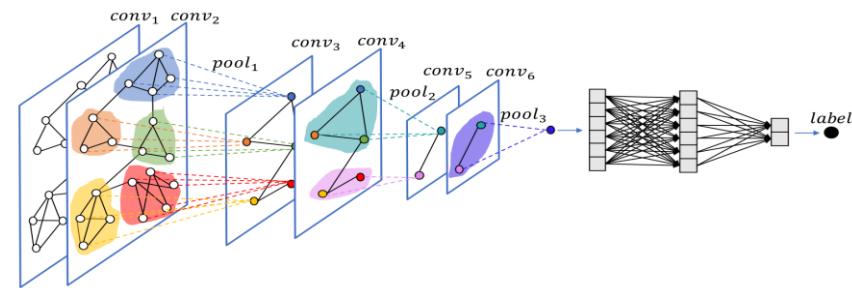


## Differentiable Graph Pooling (DIFFPOOL)[2]



Learn the cluster assignment matrix to aggregate the node representations in a hierarchical way.

## EigenPooling [3]



Incorporate the node features and local structures to obtain a better assignment matrix.

# GNN Implementation: Message Passing Framework

- Message Passing Framework:
  - **Step 1:** Gather and transform the messages from neighbors:

$$\mathbf{m}_i^{(l+1)} = \text{AGG}(\{M^{(l+1)}(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, \mathbf{e}_{i,j}) \mid j \in N(v_i)\})$$

The message generation function.  
**Input:** the state of current node, the state of the neighbor node and the edge features.

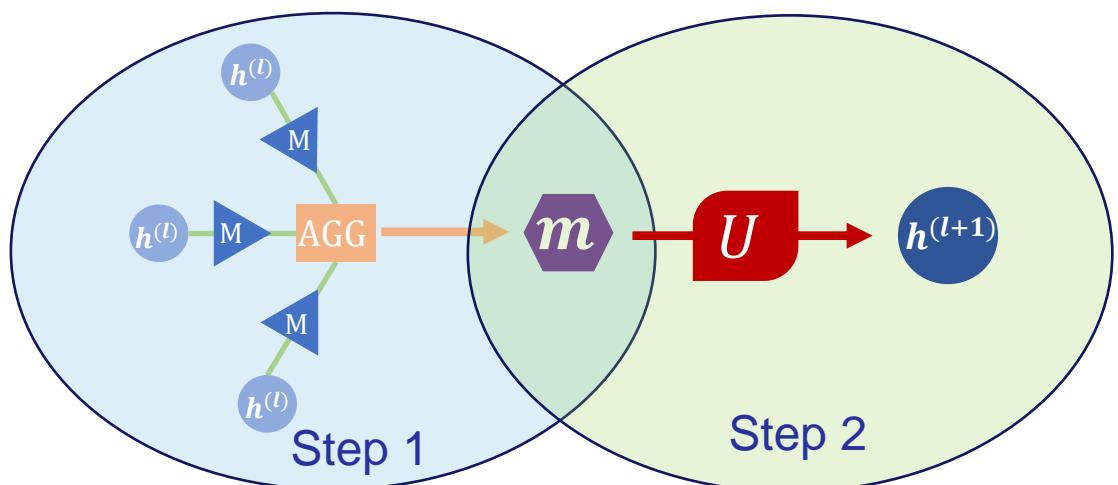
The neighborhood set of node.  
E.g. 1-hop neighbors.

The aggregation function.  
E.g. SUM/MEAN/LSTM

- **Step 2:** Update the state of the target node.

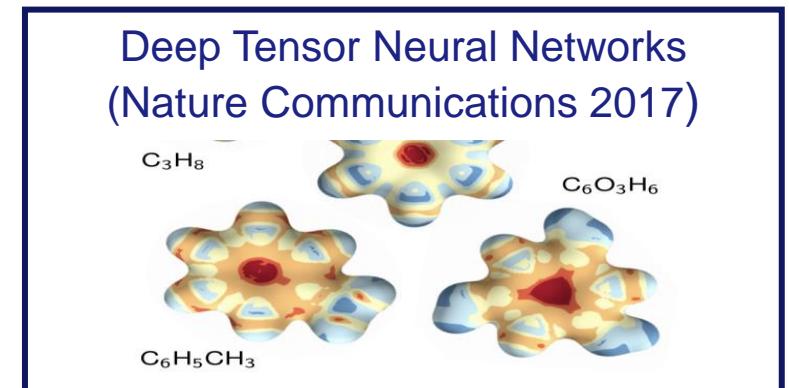
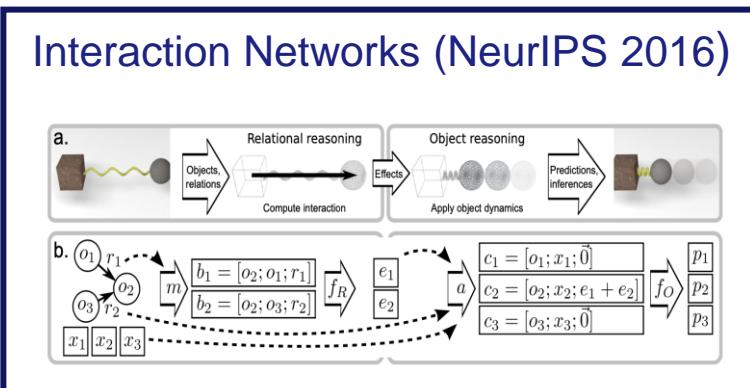
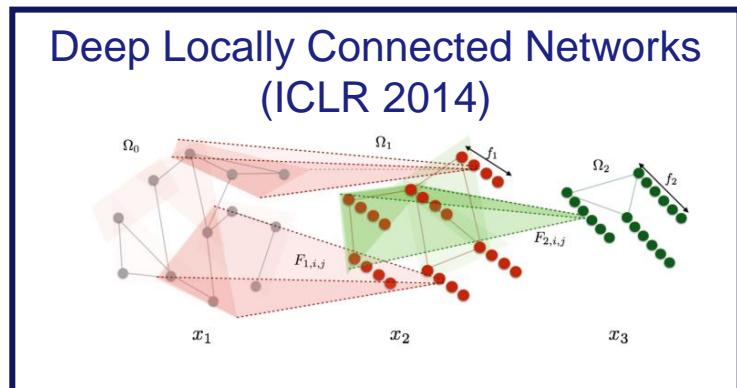
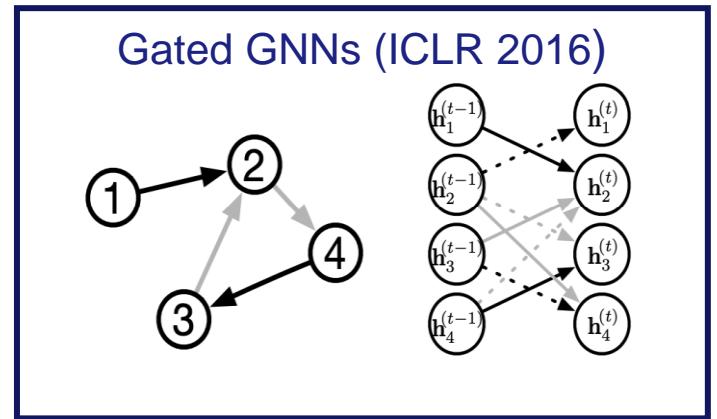
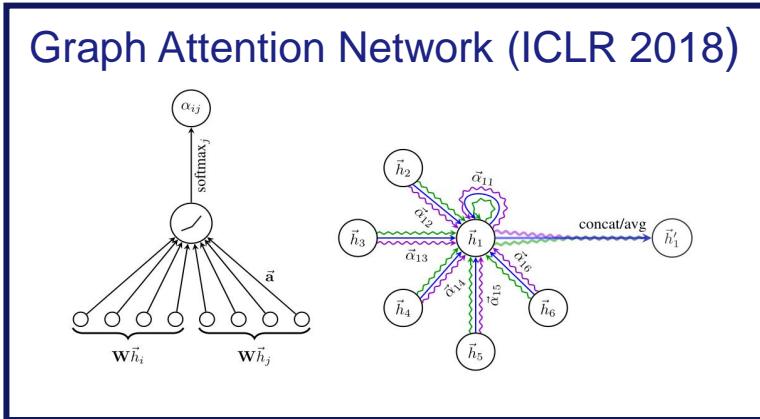
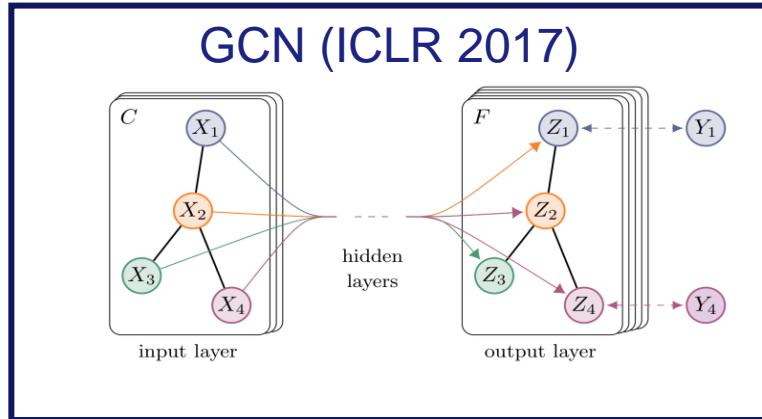
$$\mathbf{h}_i^{(l+1)} = U^{(l+1)}(\mathbf{h}_i^{(l)}, \mathbf{m}_i^{(l+1)})$$

The state update function.



# Realizations of Message Passing Neural Nets

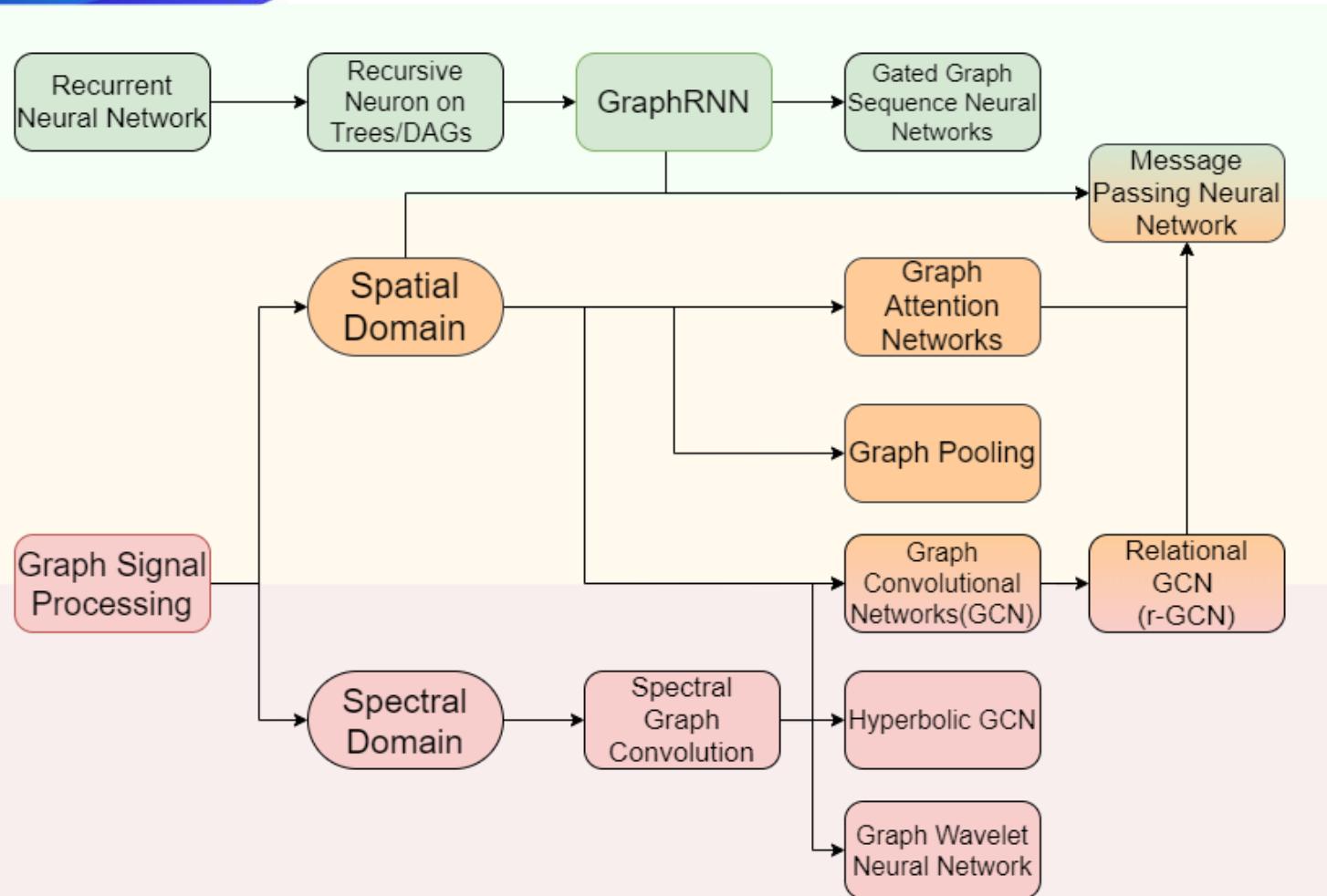
Most of current GNNs can be formulated as a message passing process.



Wang, Minjie, et al. "Deep graph library: A graph-centric, highly-performant package for graph neural networks." *arXiv preprint arXiv:1909.01315* (2019).

Fey, Matthias and Lenssen, Jan Eric Fast Graph Representation Learning with PyTorch Geometric. (2019). , cite arxiv:1903.02428.

# Summary



## Advanced topics

Training Deep GNNs

Scalability of GNNs

Robustness of GNNs

Self/Un-Supervised Learning of GNNs

....



Tencent  
AI Lab



# Training Deep GNNs

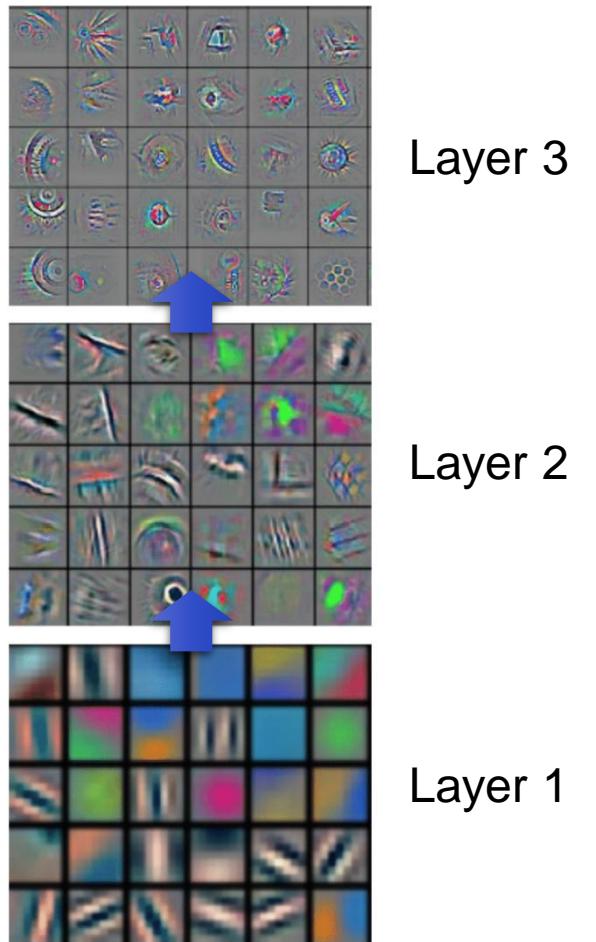
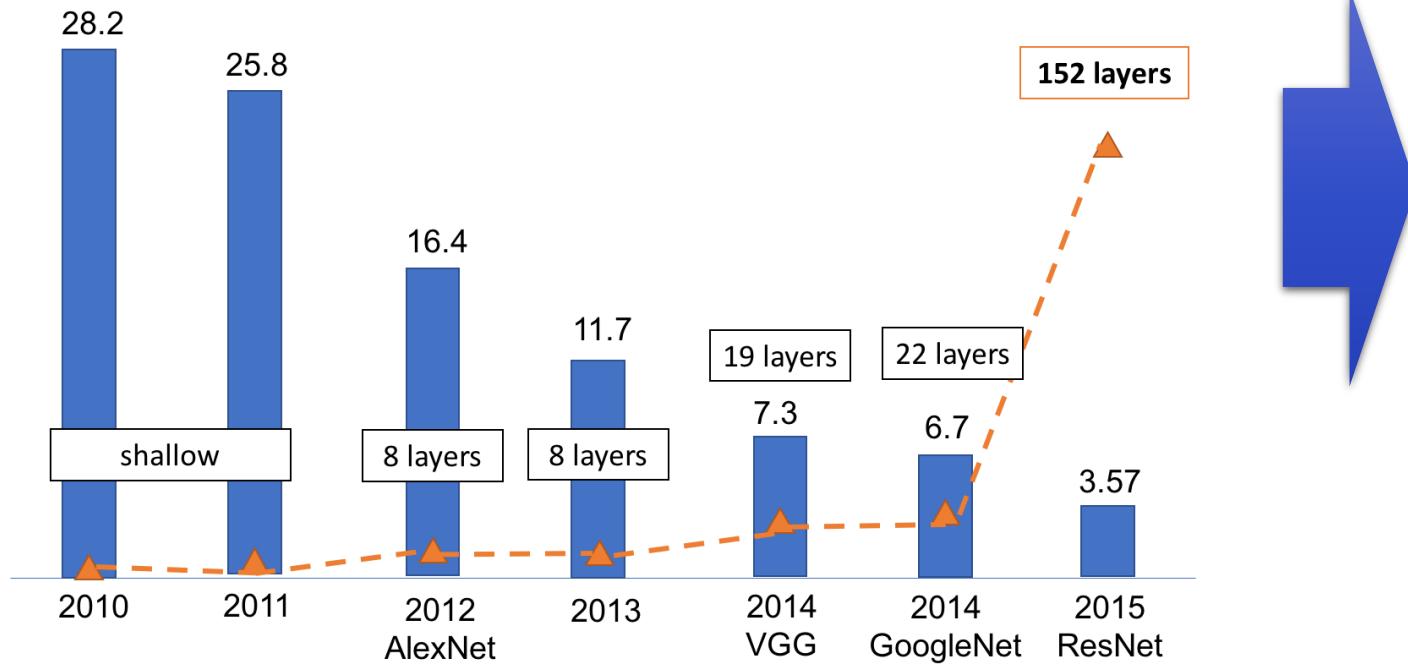


# Training Deep GNNs

- Why do we need deep GNNs?
- Can GNNs simply go deep?
- What impedes GNNs to go deep?
- How to make GNNs deep?

# The Power of Deep DNNs

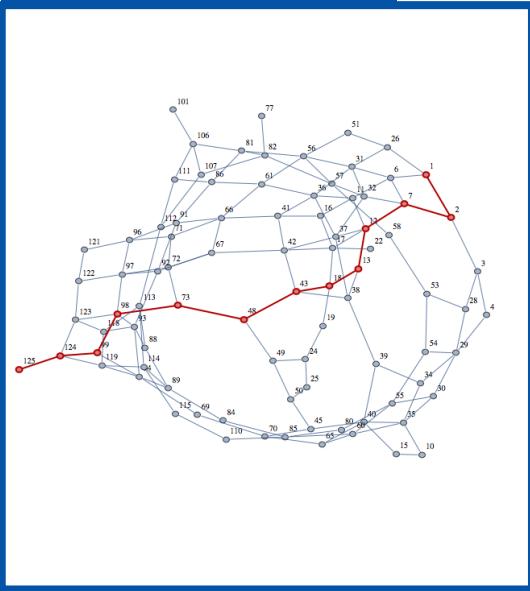
- Unprecedented success of deep DNNs in computer vision
- Deep DNNs enable larger receptive fields
- Deep DNNs enable more expressivity



# The Power of Deep GNNs

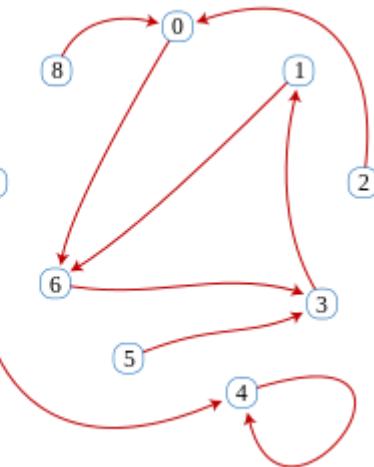
- Do GNNs need deep structures to enable larger receptive fields, too?
  - What limits the expressive power of GNNs?
    - The depth  $d$
    - The width  $w$
  - GNNs significantly lose their power when  $capacity, dw$ , is restricted

# Shortest Path

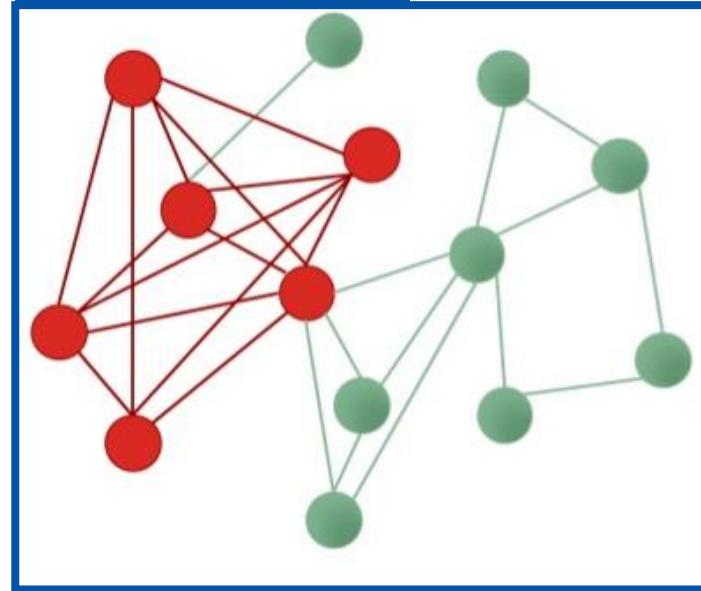


## Cycle Detection

$x$	$f(x)$
0	6
1	6
2	0
3	1
4	4
5	3
6	3
7	4
8	0



## Subgraph

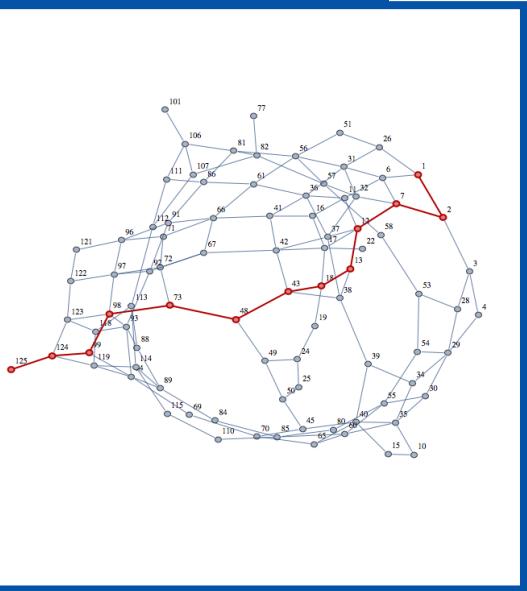


# The Power of Deep GNNs

- ◀ Do GNNs need deep structures to enable larger receptive fields, too?
- ◀ What limits the expressive power of GNNs?
  - ◀ The depth  $d$
  - ◀ The width  $w$
- ◀ GNNs significantly lose their power when *capacity*,  $dw$ , is restricted

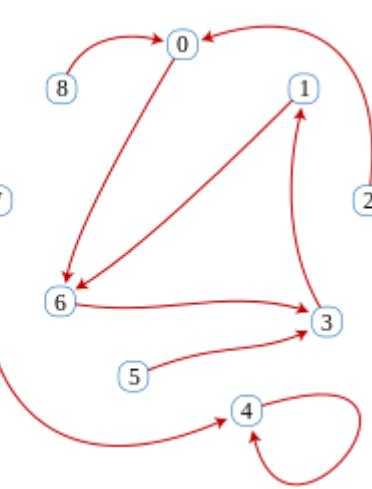
**Yes**

Shortest Path

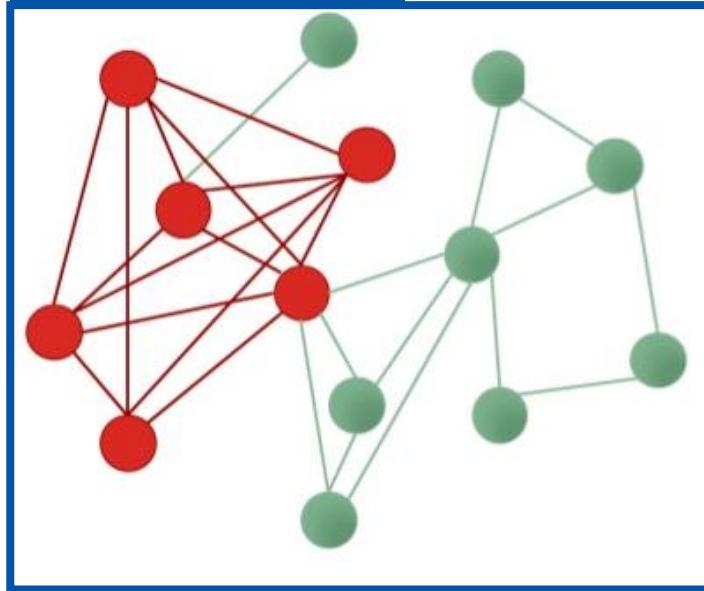


Cycle Detection

$x$	$f(x)$
0	6
1	6
2	0
3	1
4	4
5	3
6	3
7	4
8	0



Subgraph



# The Power of Deep GNNs

## ► The boundary of capacity for different problems

<i>problem</i>	<i>bound</i>	<i>problem</i>	<i>bound</i>
cycle detection (odd)	$d_w = \Omega(n/\log n)$	shortest path	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$
cycle detection (even)	$d_w = \Omega(\sqrt{n}/\log n)$	max. indep. set	$d_w = \Omega(n^2/\log^2 n)$ for $w = O(1)$
subgraph verification*	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	min. vertex cover	$d_w = \Omega(n^2/\log^2 n)$ for $w = O(1)$
min. spanning tree	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	perfect coloring	$d_w = \Omega(n^2/\log^2 n)$ for $w = O(1)$
min. cut	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	girth 2-approx.	$d_w = \Omega(\sqrt{n}/\log n)$
diam. computation	$d_w = \Omega(n/\log n)$	diam. $3/2$ -approx.	$d_w = \Omega(\sqrt{n}/\log n)$

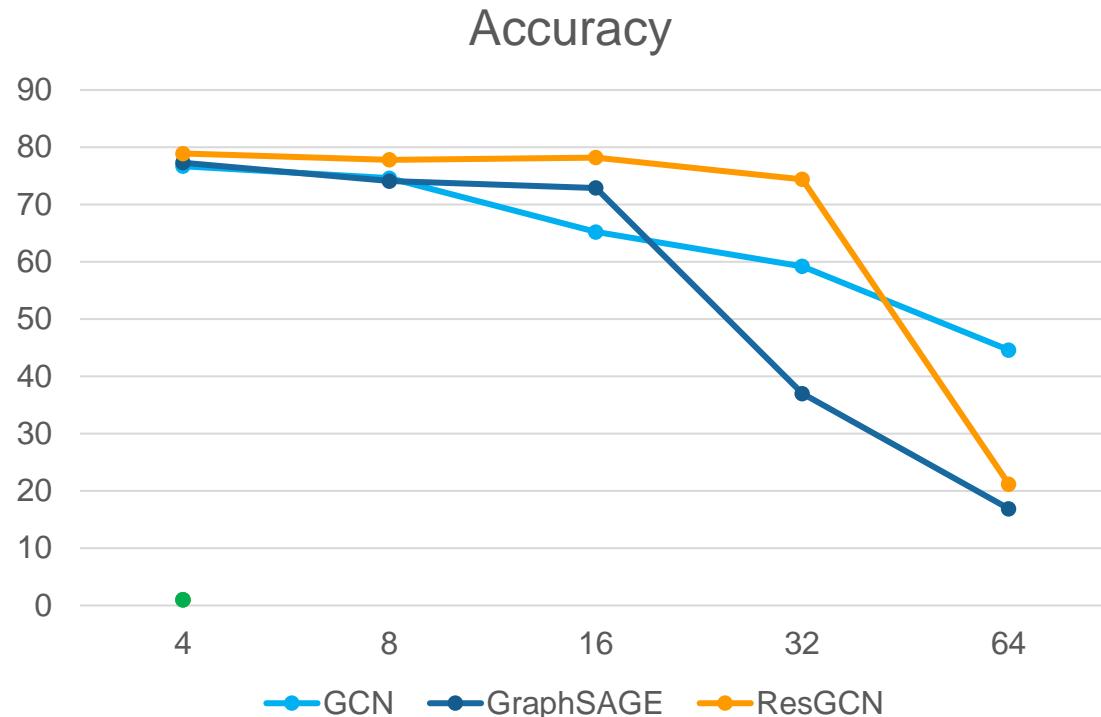
(Loukas, ICLR'20)

# Training Deep GNNs

- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
  - ◀ **GCN**: Basic GCN
  - ◀ **GraphSAGE**: GCN with improved **aggregation**
  - ◀ **ResGCN**: leverage idea from **ResNet**
- ◀ What impedes GNNs to go deep?
- ◀ How to make GNNs deep?

# GNNs are Shallow

But can they really go deep? Not all



Citeseer	4 layers	16 layers	64 layers
GCN	76.7	65.2	44.6
GraphSAGE	77.3	72.9	16.9
ResGCN	78.9	78.2	21.2

(Rong et al, ICLR'20)

# Training deep GNNs

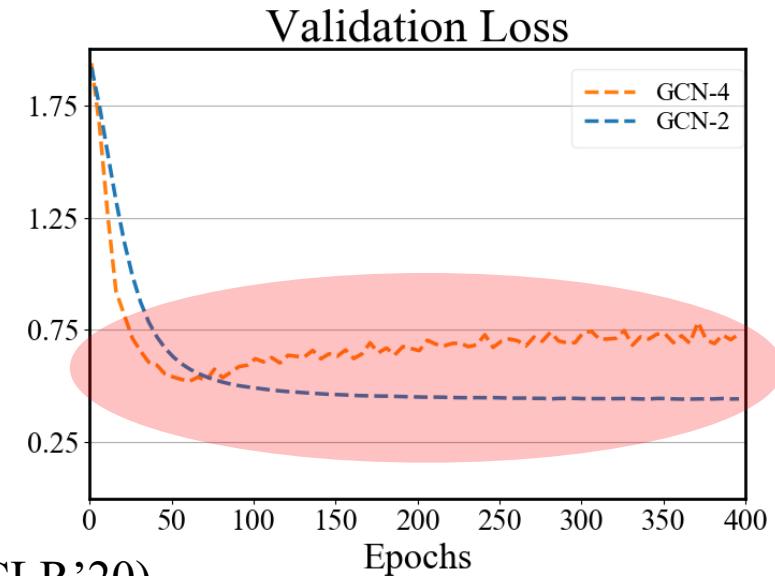
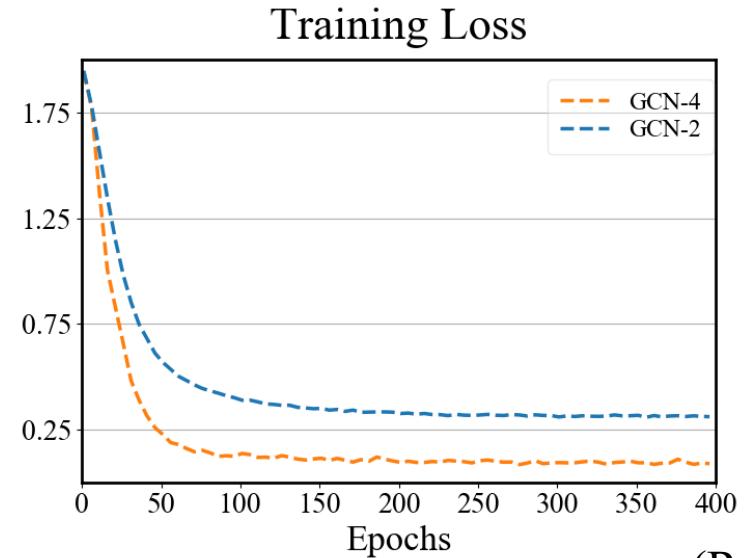
- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ Overfitting (Common)
  - ◀ Training dynamics (Common)
  - ◀ Over-smoothing (Graph Specific)
- ◀ How to make GNNs deep?

# Training deep GNNs

- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ **Overfitting (Common)**
  - ◀ Training dynamics (Common)
  - ◀ Over-smoothing (Graph Specific)
- ◀ How to make GNNs deep?

# Overfitting

## ◆ GNNs suffer from Overfitting



(Rong et al, ICLR'20)

## ◆ Too many parameters are established but only few of data points are provided

$$O(dh^2)$$

# Training deep GNNs

- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ Overfitting (Common)
  - ◀ **Training dynamics (Common)**
  - ◀ Over-smoothing (Graph Specific)
- ◀ How to make GNNs deep?

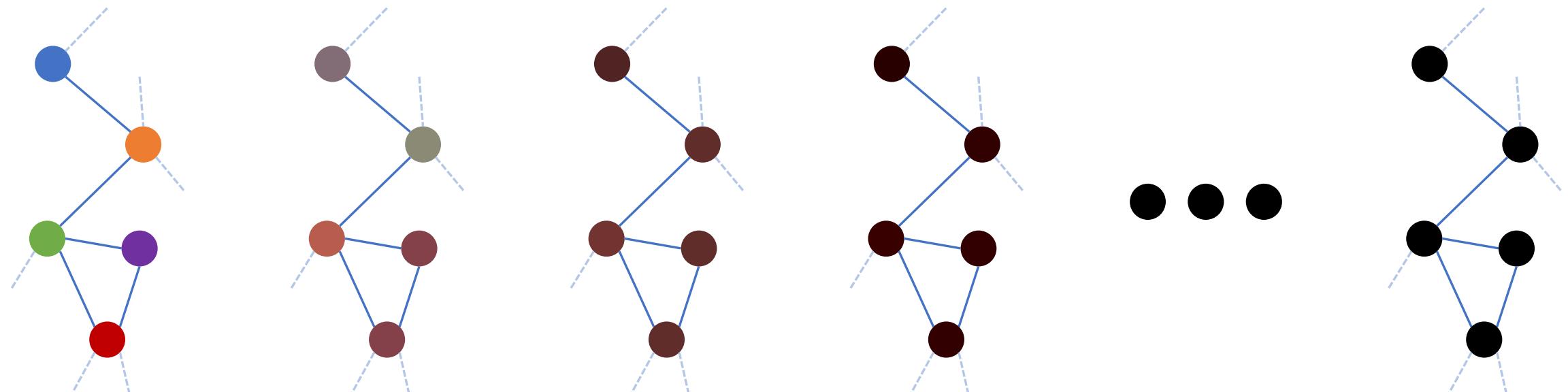


# Training dynamics

$l$ -layers gradient

$$\frac{dH_{l+1}}{dH_l} \cdot \frac{dH_l}{dH_{l-1}} \cdot \dots \cdot \frac{dH_0}{dW_0} \leq (s_l \lambda_{m+1}) \cdot (s_{l-1} \lambda_{m+1}) \cdot \dots \cdot \frac{dH_0}{dW_0}$$

The gradients vanish as the model go deep because  $s_1 \dots s_l \lambda_{m+1} < 1$



RGB as Features

Layer 1

Layer 100

Layer 200

Layer 500  
RGB=[0,0,0]

# Training deep GNNs

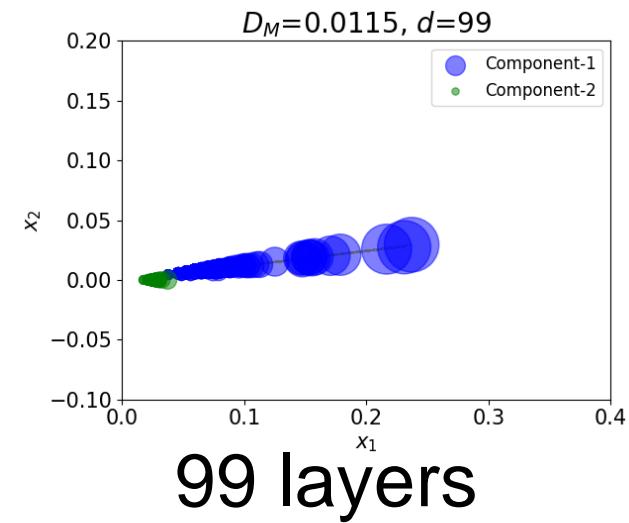
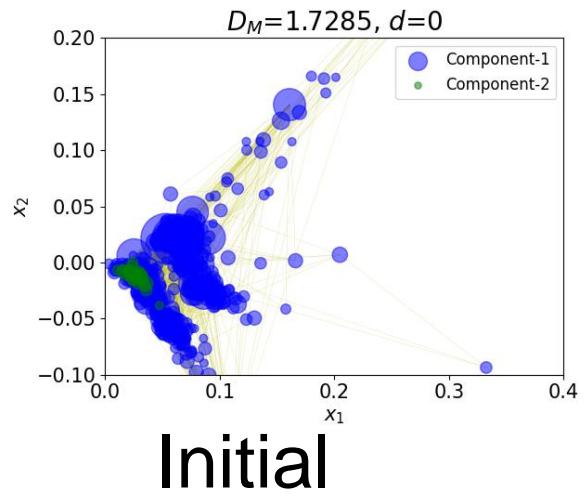
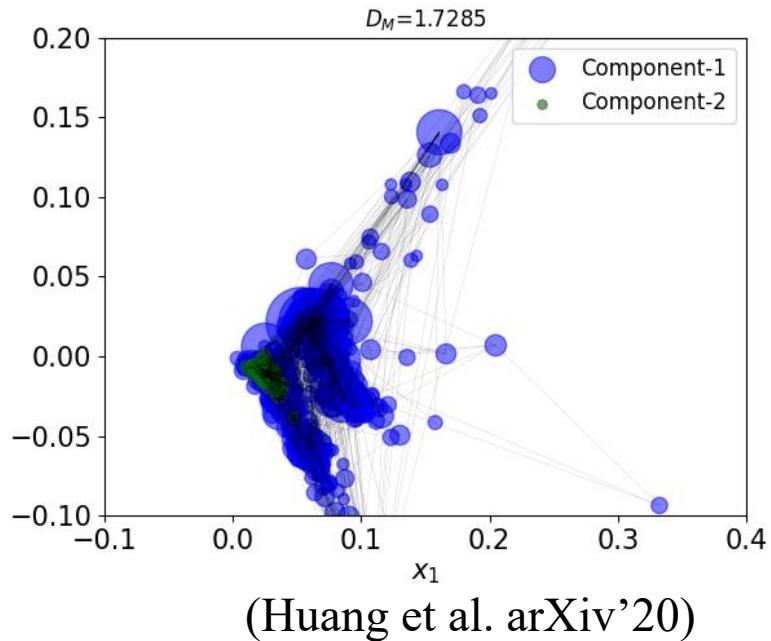
- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ Overfitting (Common)
  - ◀ Training dynamics (Common)
  - ◀ **Over-smoothing (Graph Specific)**
- ◀ How to make GNNs deep?

# Over-Smoothing

► GNNs suffers from over-smoothing

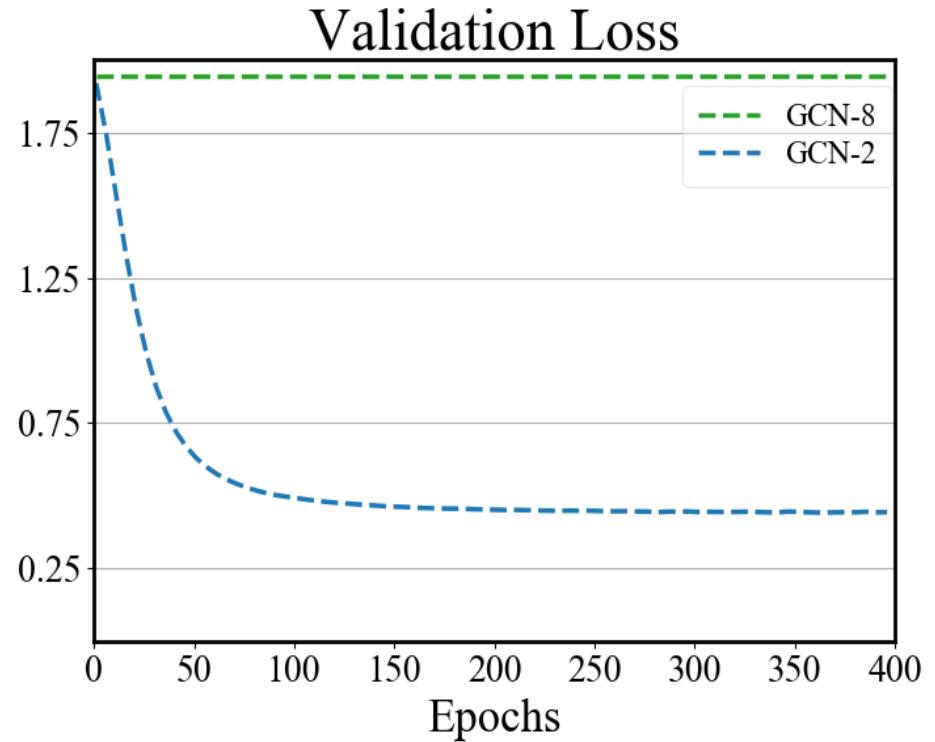
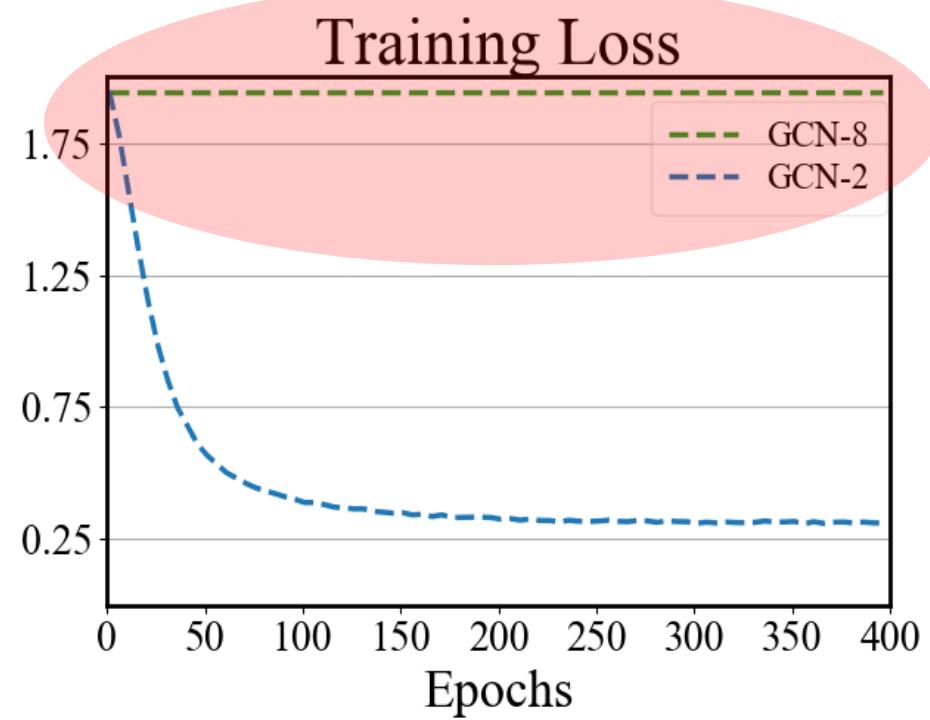
► Over-smoothing: node representations become **less distinguishable** with each other when the depth increases

(Li et al. AAAI'18; Chen et al. AAAI'20; Oono et al. ICLR'20; Rong et al. ICLR'20; Huang et al. arXiv'20)



# Over-Smoothing

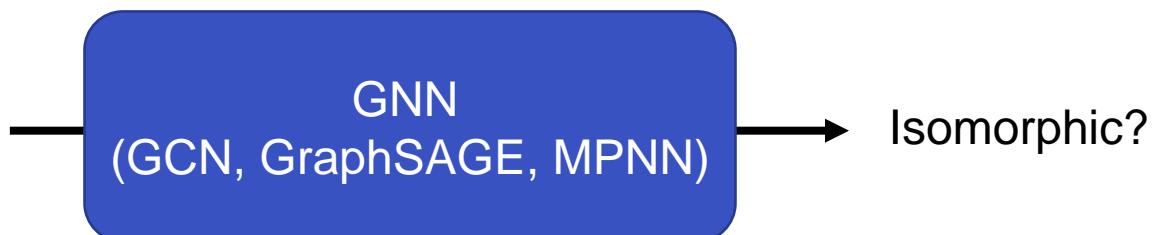
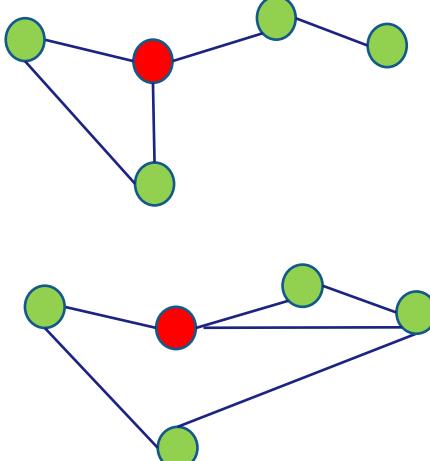
Over-smoothing also impedes training.



# Over-Smoothing

## Why GNN works?

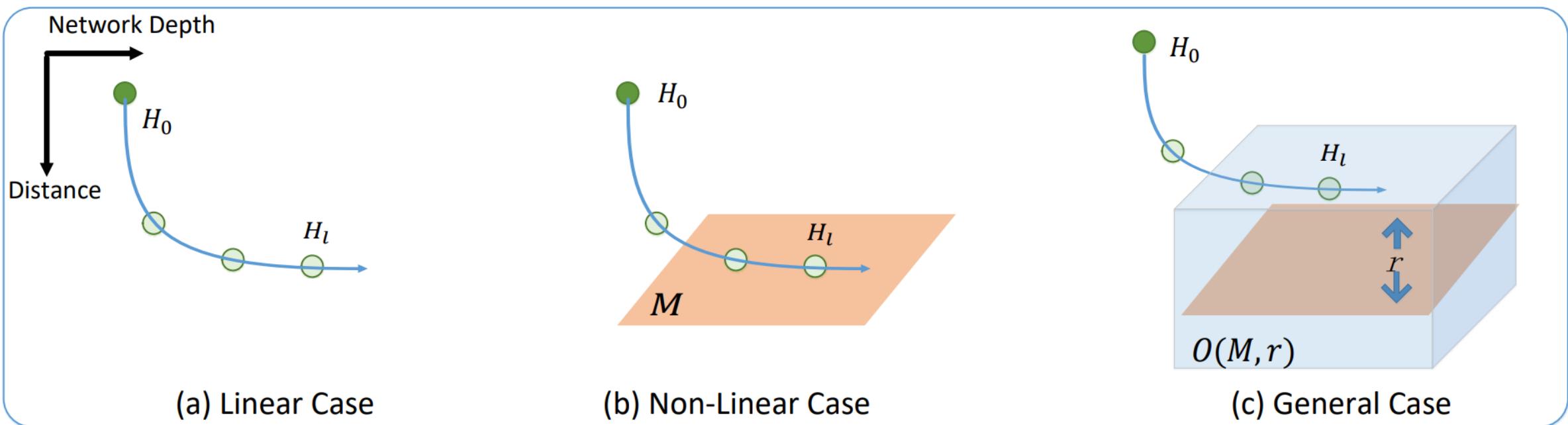
- By message passing, GNN is able to capture the local structure;
- Several works (Xu et al., 2019, Murphy et al., 2019) show that GNN is equivalent to the Weisfeiler-Lehman (WL) test under a careful design



# Over-Smoothing

## When GCNs fail?

- With linear activation
- With ReLU activation
- With ReLU and bias



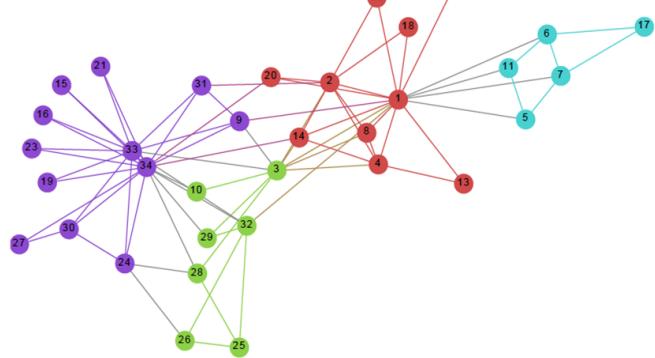
# Over-Smoothing of Linear GCN

- When GCNs fail?
- With linear activation

$l$ -step Random Walk

$$y = \hat{A}^l p_0 \text{ where } p_0(i) = \boxed{1/d(i)}$$

Probability of walking



## Random Walks on Graph

- $V_{26} - V_{25} - V_{32} - V_3 - V_{10} \dots$
- $V_5 - V_7 - V_{17} - V_6 - V_{11} \dots$
- $V_{31} - V_{33} - V_{21} - V_{33} - V_{15}$

# Over-Smoothing of Linear GCN

- When GCNs fail?
  - With linear activation

$l$ -step Random Walk

$$y = \hat{A}^l p_0 \text{ where } p_0(i) = 1/d(i)$$

$l$ -layer GCNs

$$Y = \hat{A}^l [XW]$$

Learnable Probability

# Over-Smoothing of Linear GCN

- 💡 When GCNs fail?
- 💡 With linear activation

*l*-step Random Walk

$$y = \hat{A}^l p_0 \text{ where } p_0(i) = 1/d(i)$$

*l*-layer GCNs

$$Y = \hat{A}^l X W$$

Eigen decomposition

$$Y = \sum_{i=1}^n (\lambda_i u_i u_i^\top)^l X W$$



# Over-Smoothing of Linear GCN

## Rewrite eigen decomposition

Eigen decomposition

$$(\lambda_1 u_1 u_1^\top)^l X W + \cdots (\lambda_m u_m u_m^\top)^l X W + (\lambda_{m+1} u_{m+1} u_{m+1}^\top)^l X W + \cdots (\lambda_n u_n u_n^\top)^l X W$$

# Over-Smoothing of Linear GCN

💡 Rewrite eigen decomposition

Eigen decomposition

$$(\lambda_1 u_1 u_1^\top)^l X W + \cdots (\lambda_m u_m u_m^\top)^l X W + (\lambda_{m+1} u_{m+1} u_{m+1}^\top)^l X W + \cdots (\lambda_n u_n u_n^\top)^l X W$$

💡 Suppose graph  $\mathcal{G}$  has  $m$  connected components. It indicates

Eigenvalues

$$1 = \lambda_1 = \cdots = \lambda_m > \lambda_{m+1} > \cdots > \lambda_n > -1$$

# Over-Smoothing of Linear GCN

💡 Rewrite eigen decomposition

Eigen decomposition

$$(\lambda_1 u_1 u_1^\top)^l X W + \dots + (\lambda_m u_m u_m^\top)^l X W + (\lambda_{m+1} u_{m+1} u_{m+1}^\top)^l X W + \dots + (\lambda_n u_n u_n^\top)^l X W$$

💡 Suppose graph  $\mathcal{G}$  has  $m$  connected components. It indicates

Eigenvalues

$$1 = \lambda_1 = \dots = \lambda_m > \lambda_{m+1} > \dots > \lambda_n > -1$$

💡 When  $l \rightarrow +\infty$ ,  $\lambda_{m+1}, \dots, \lambda_n \rightarrow 0$

Convergence

$$Y = \lim_{l \rightarrow \infty} [\lambda_1^l u_1 u_1^\top X W + \dots + \lambda_m^l u_m u_m^\top X W + \lambda_{m+1}^l u_{m+1} u_{m+1}^\top X W + \dots + \lambda_n^l u_n u_n^\top X W]$$


 1


 1


 0


 0

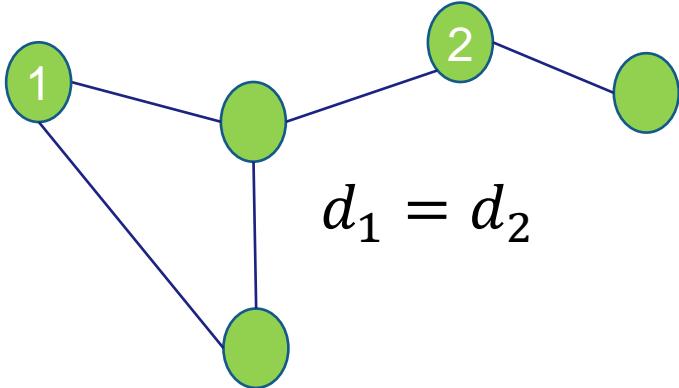
# Over-Smoothing of Linear GCN

When  $l \rightarrow +\infty$ ,  $\lambda_{m+1}, \dots, \lambda_n \rightarrow 0$

## Convergence

$$Y = u_1(u_1^\top XW) + \dots + u_m(u_m^\top XW), \text{ where } u_i(j) = \frac{1}{d_j^2} \cdot \delta(\text{node } j \text{ in component } i)$$

Node 1 is **indistinguishable** with node 2



	degree	$Y$	
Node 1	$d_1^{\frac{1}{2}}z_1$	$d_1^{\frac{1}{2}}z_2$	$d_1^{\frac{1}{2}}z_3$
Node 2	$d_2^{\frac{1}{2}}z_1$	$d_2^{\frac{1}{2}}z_2$	$d_2^{\frac{1}{2}}z_3$
			...

The nodes within the same connected component are distinguishable only by their degrees

# Over-Smoothing of Non-Linear GCN

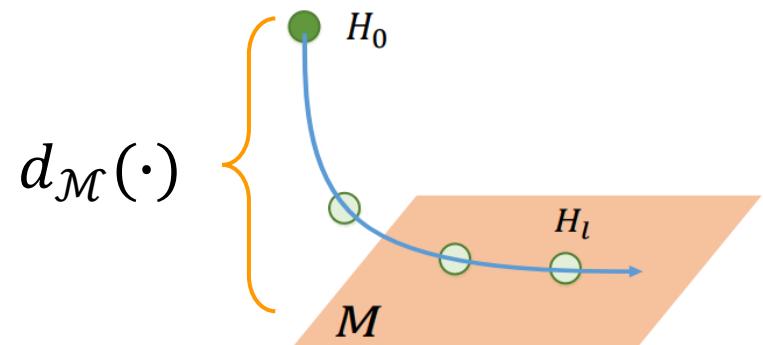
## With ReLU activation

- Similar to the linear case, but hard to derive the exact convergence point. Require the notion of subspace (Oono et al., ICLR'20):

$\mathcal{M}$  subspcae

**Definition 1** (subspace). Let  $\mathcal{M} := \{\mathbf{E}\mathbf{C} | \mathbf{C} \in \mathbb{R}^{M \times C}\}$  be an  $M$ -dimensional subspace in  $\mathbb{R}^{N \times C}$ , where  $\mathbf{E} \in \mathbb{R}^{N \times M}$  is orthogonal, i.e.  $\mathbf{E}^T \mathbf{E} = \mathbf{I}_M$ , and  $M \leq N$ .

It is proved that (Oono et al., ICLR'20):  
an infinite-layer GCN will converge to a certain point within a subspace  $\mathcal{M}$



(b) Non-Linear Case



# Over-Smoothing of Non-Linear GCN

Convergence of  $\hat{A}$

$$d_{\mathcal{M}}(\hat{A}X) \leq \lambda_{m+1} d_{\mathcal{M}}(X), \lambda_{m+1} < 1$$



Convergence

$$d_{\mathcal{M}}(H_{l+1}) = d_{\mathcal{M}}(\sigma(\hat{A}H_lW)) \leq \lambda_{m+1} s d_{\mathcal{M}}(H_l)$$

# Over-Smoothing of Non-Linear GCN

Convergence of  $\hat{A}$

$$d_{\mathcal{M}}(\hat{A}X) \leq \lambda_{m+1} d_{\mathcal{M}}(X), \lambda_{m+1} < 1$$

Convergence of  $W$

$$d_{\mathcal{M}}(XW_l) \leq s d_{\mathcal{M}}(X), s \leq 1$$

Convergence

$$d_{\mathcal{M}}(H_{l+1}) = d_{\mathcal{M}}(\sigma(\hat{A}H_l W)) \leq \lambda_{m+1} s d_{\mathcal{M}}(H_l)$$



# Over-Smoothing of Non-Linear GCN

Convergence of  $\hat{A}$

$$d_{\mathcal{M}}(\hat{A}X) \leq \lambda_{m+1} d_{\mathcal{M}}(X), \lambda_{m+1} < 1$$

Convergence of  $W$

$$d_{\mathcal{M}}(XW) \leq s d_{\mathcal{M}}(X), \quad s_l \leq 1$$

Convergence of ReLU

$$d_{\mathcal{M}}(\sigma(X)) \leq d_{\mathcal{M}}(X)$$

Convergence

$$d_{\mathcal{M}}(H_{l+1}) = d_{\mathcal{M}}(\sigma(\hat{A}H_l W)) \leq \lambda_{m+1} s d_{\mathcal{M}}(H_l)$$



# Over-Smoothing of GCNs with bias

- With ReLU and bias

GCNs with bias

$$H_{l+1} = \sigma(\hat{A}H_lW_l + b_l)$$

# Over-Smoothing of GCNs with bias

- $H_L$  converges to a certain sub-cube  $\mathcal{O}(\mathcal{M}, r)$  with ReLU and bias

GCNs with bias

$$H_{l+1} = \sigma(\hat{A}H_lW + b)$$

Convergence of bias

$$d_{\mathcal{M}}(H_{l+1}) \leq \lambda_{m+1} s d_{\mathcal{M}}(H_l) + d_{\mathcal{M}}(b)$$

# Over-Smoothing of GCNs with bias

- With ReLU and bias

GCNs with bias

$$H_{l+1} = \sigma(\hat{A}H_lW + b)$$

Convergence of bias

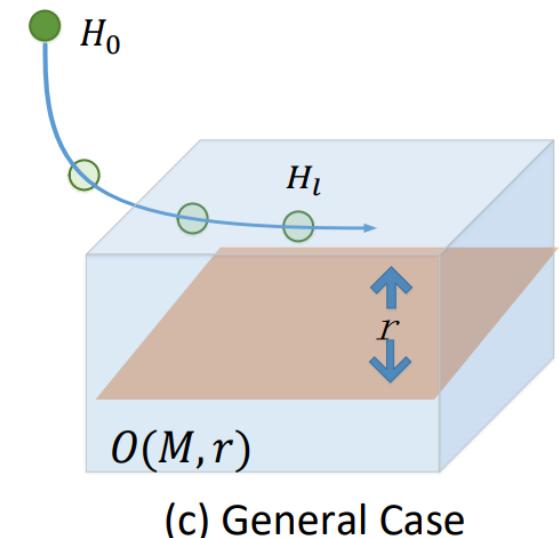
$$d_{\mathcal{M}}(H_{l+1}) \leq \lambda_{m+1}s d_{\mathcal{M}}(H_l) + d_{\mathcal{M}}(b)$$

GCN with bias

$H_l$  converges to a certain sub-cube:

$$\mathcal{O}(\mathcal{M}, r) = \{H_l | d_{\mathcal{M}}(H_l) < \frac{d_{\mathcal{M}}(b)}{1 - \lambda_{m+1}s}\}$$

(Huang et al. arXiv'20)



# Summary

Linear GCN

$H_l$  converges to a certain point that can be exactly derived

Non-linear GCN

$H_l$  converges to a certain point within a certain subspace

GCN with bias

$H_l$  converges to a certain point within a certain sub-cube

# Other methods to measure over-smoothing

- One can explicitly measure over-smoothing using distance between node pairs (Chen et al., 2020)

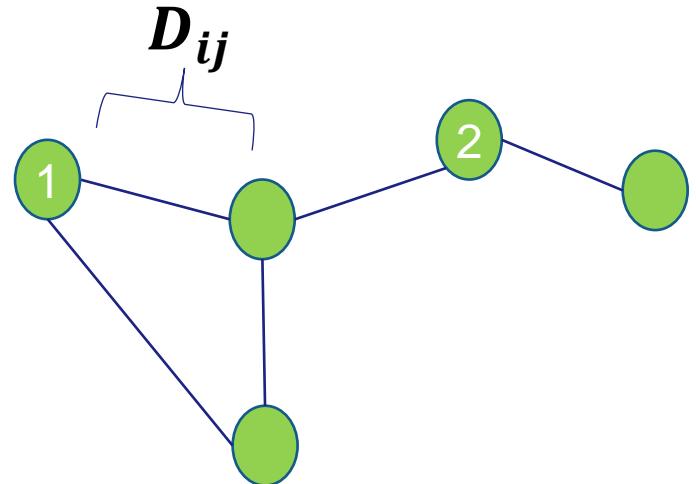
Pair Distance

$$D_{ij} = 1 - \frac{H_i^T H_j}{\|H_i\| \|H_j\|}$$

The i-th row  
of the output

Mean Average Distance

$$\text{MAD}_i = \frac{\sum_{i=0}^n \bar{D}_i^{tgt}}{\sum_{i=0}^n \mathbb{I}(\bar{D}_i^{tgt})}$$



- Other metrics see PairNorm (Zhao et al., 2020); GroupNorm (Zhou et al., 2020)

Chen, et al. "Measuring and Relieving the Over-smoothing Problem for Graph Neural Networks from the Topological View." AAAI 2020

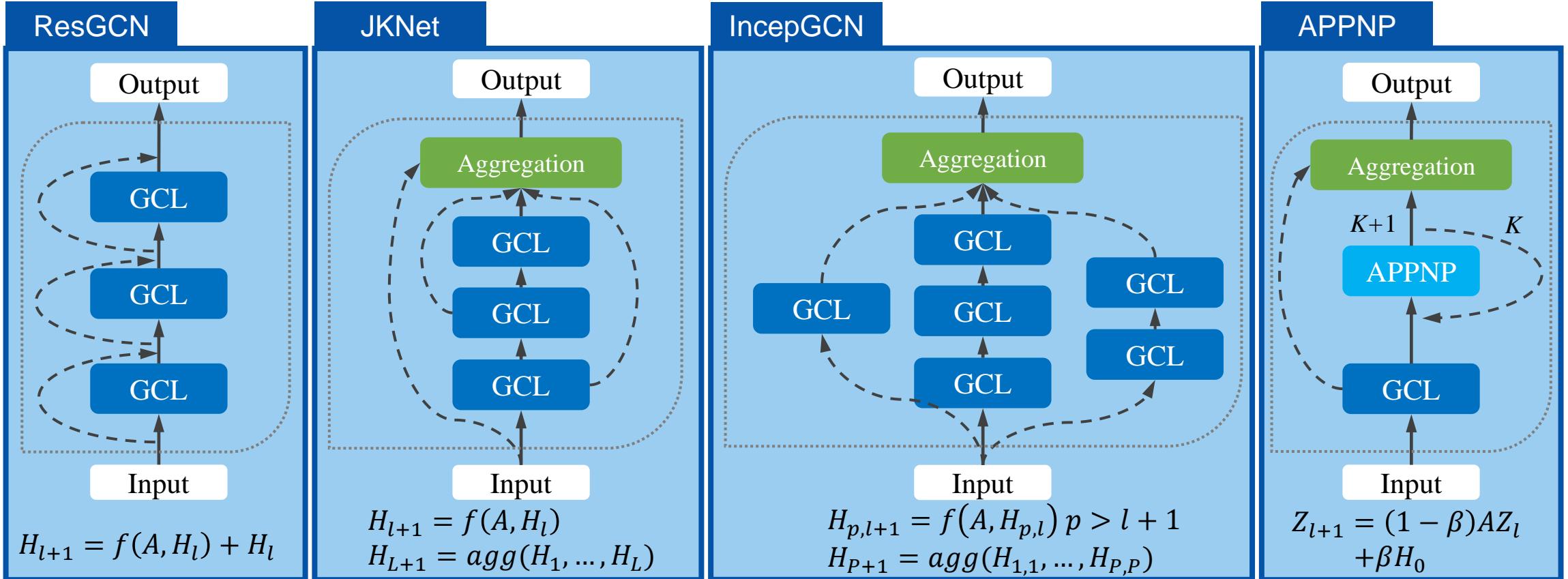
Zhao, Lingxiao, and Leman Akoglu. "PairNorm: Tackling Oversmoothing in GNNs." ICLR. 2020.

Zhou et al. "Towards Deeper Graph Neural Networks with Differentiable Group Normalization Kaixiong Zhou Texas A&M University zkxiong@tamu.edu Xiao Huang The Hong." NIPS 2020.

# Training deep GNNs

- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ Overfitting (Common)
  - ◀ Training dynamics (Common)
  - ◀ Over-smoothing (Graph Specific)
- ◀ How to make GNNs deep?
  - ◀ Architecture refinement
  - ◀ Manipulating input (DropEdge)
  - ◀ Layer normalizations

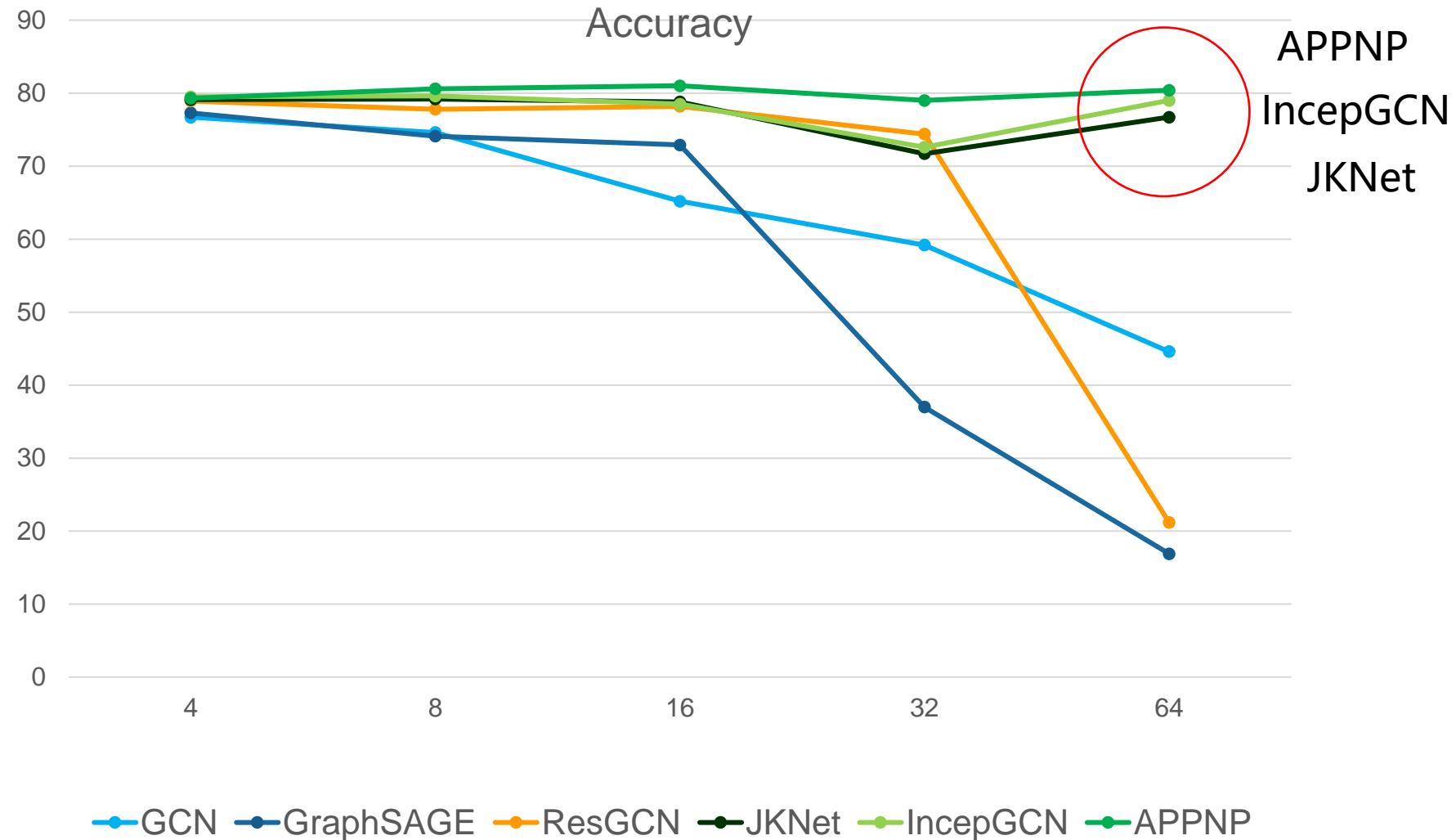
# Shortcuts in Structures



Other architectures including SGC (Wu et al., 2019), GCNII (Chen et al., 2020), etc.



# Shortcuts in Structures



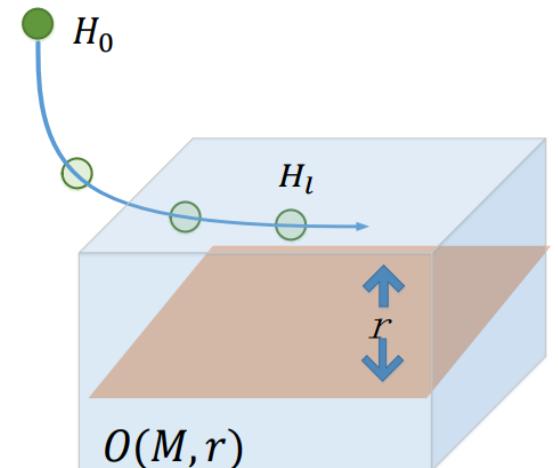
# Shortcuts in Structures

- Residual connections are helpful (akin to CNNs)
- Do residual connections alleviate over-smoothing?

GCN with bias

$$d_{\mathcal{M}}(H_{l+1}) - r \leq \nu(d_{\mathcal{M}}(H_l) - r)$$

Indeed, general GCNs converge to a certain sub-cube  $O(\mathcal{M}, r)$  with speed  $\nu^{-1}$  and radius  $r$ .



(c) General Case

# Shortcuts in Structures

General Case

Converging to sub-cube with speed  $\nu^{-1}$  and radius  $r$

Basic GCN

Generic GCN

$$\nu = s\lambda_{m+1}$$

$$r = 0$$

GCN with bias

$$\nu = s\lambda_{m+1}$$

$$r = \frac{d_{\mathcal{M}}(b)}{1 - \nu}$$

Different Structures

ResGCN

$$\nu = s\lambda_{m+1} + \alpha$$

$$r = 0$$

APPNP

$$\nu = (1 - \beta) \lambda_{m+1}$$

$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - \nu}$$

# Shortcuts in Structures

General Case

General models converge to  $O(\mathcal{M}, r)$  with speed  $\nu^{-1}$ , when  $l \rightarrow \infty$

Basic GCN

Generic GCN

$$\nu = s\lambda_{m+1}$$

$$r = 0$$

GCN with bias

$$\nu = s\lambda_{m+1} \quad r \text{ is nonzero}$$

$$r = \frac{d_{\mathcal{M}}(\mathbf{b})}{1 - \nu}$$

Different Structures

ResGCN

$$\nu = s\lambda_{m+1} + \alpha$$

$$r = 0$$

APPNP

$$\nu = (1 - \beta) \lambda_{m+1}$$

$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - \nu}$$

# Shortcuts in Structures

General Case

General models converge to  $O(\mathcal{M}, r)$  with speed  $\nu^{-1}$ , when  $l \rightarrow \infty$

Basic GCN

Generic GCN

$$\nu = s\lambda_{m+1}$$

$$r = 0$$

GCN with bias

$$\nu = s\lambda_{m+1}$$

$$r = \frac{d_{\mathcal{M}}(b)}{1 - \nu}$$

Different Structures

ResGCN

$\nu$  is increased

$$\nu = s\lambda_{m+1} + \alpha$$

$$r = 0$$

APPNP

$$\nu = (1 - \beta) \lambda_{m+1}$$

$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - \nu}$$

# Shortcuts in Structures

General Case

General models converge to  $O(\mathcal{M}, r)$  with speed  $\nu^{-1}$ , when  $l \rightarrow \infty$

Basic GCN

Generic GCN

$$\nu = s\lambda_{m+1}$$

$$r = 0$$

GCN with bias

$$\nu = s\lambda_{m+1}$$

$$r = \frac{d_{\mathcal{M}}(b)}{1 - \nu}$$

Different Structures

ResGCN

$$\nu = s\lambda_{m+1} + \alpha$$

$$r = 0$$

APPNP

$$\nu = (1 - \beta) \lambda_{m+1}$$

$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - \nu}$$

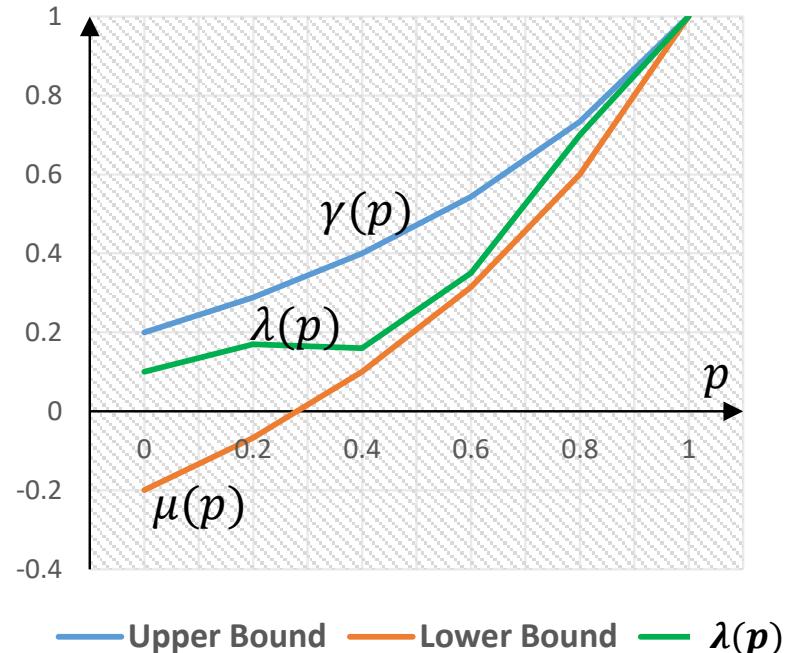
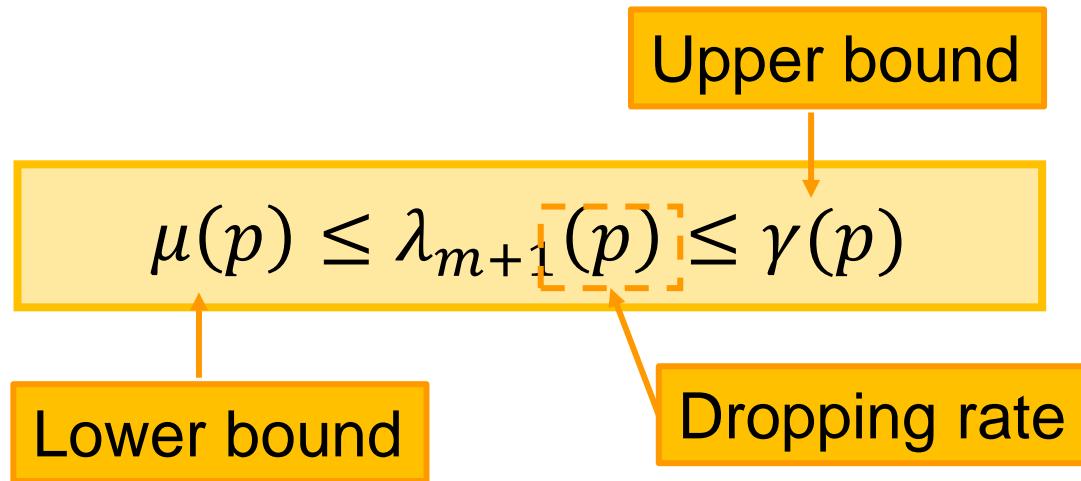
$\nu$  is decreased  
 $r$  is nonzero

# Over-smoothing Layer

- For all cases, the over-smoothing speed  $\nu^{-1}$  is controlled by  $\lambda_{m+1}$  (the second-biggest eigenvalue of normalized adjacency matrix)
- So how to increase  $\lambda_{m+1}$ ?

# Alleviate Over-Smoothing by DropEdge

- So how to increase  $\lambda_{m+1}$ ? Drop Edges!
- In expectation:



- Both  $\mu$  and  $\gamma$  monotonically increase w.r.t.  $p$ ;
- The gap  $\gamma - \mu$  monotonically decreases w.r.t.  $p$ ;

# Alleviate Over-Smoothing by DropEdge

- Huang et al., 2020 has considered the re-normalization trick in our analyses, in contrast to Rong et al., 2020

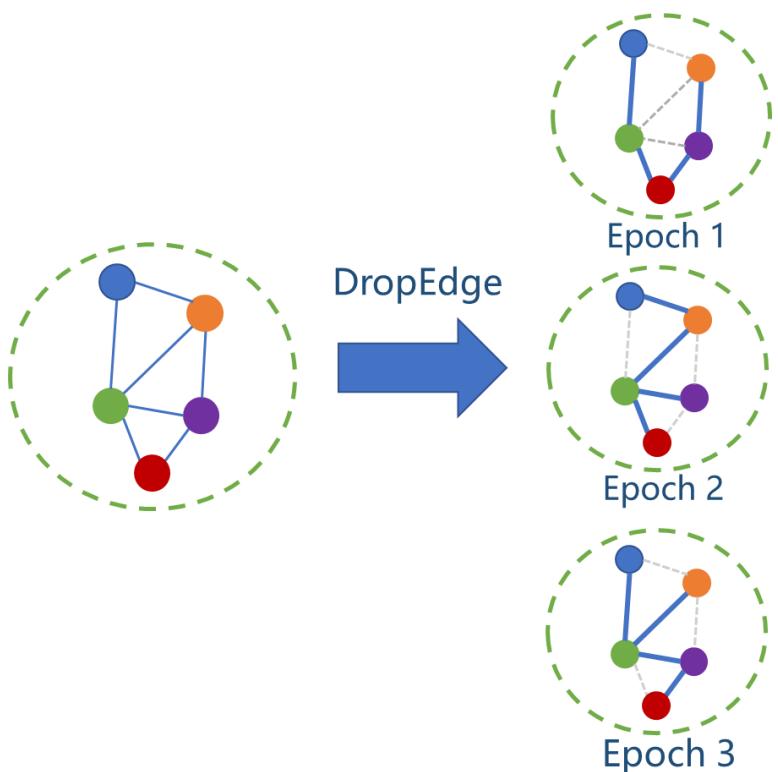
**Theorem 1.** We denote the original graph as  $\mathcal{G}$  and the one after dropping certain edges out as  $\mathcal{G}'$ . Given a small value of  $\epsilon$ , we assume  $\mathcal{G}$  and  $\mathcal{G}'$  will encounter the  $\epsilon$ -smoothing issue with regard to subspaces  $\mathcal{M}$  and  $\mathcal{M}'$ , respectively. Then, either of the following inequalities holds after sufficient edges removed.

- The relaxed smoothing layer only increases:  $\hat{l}(\mathcal{M}, \epsilon) \leq \hat{l}(\mathcal{M}', \epsilon)$ ;
- The information loss is decreased:  $N - \dim(\mathcal{M}) > N - \dim(\mathcal{M}')$ .

(Rong et al., 2020)

# Alleviate Over-Smoothing by DropEdge

Besides, DropEdge can prevent over-fitting as well!



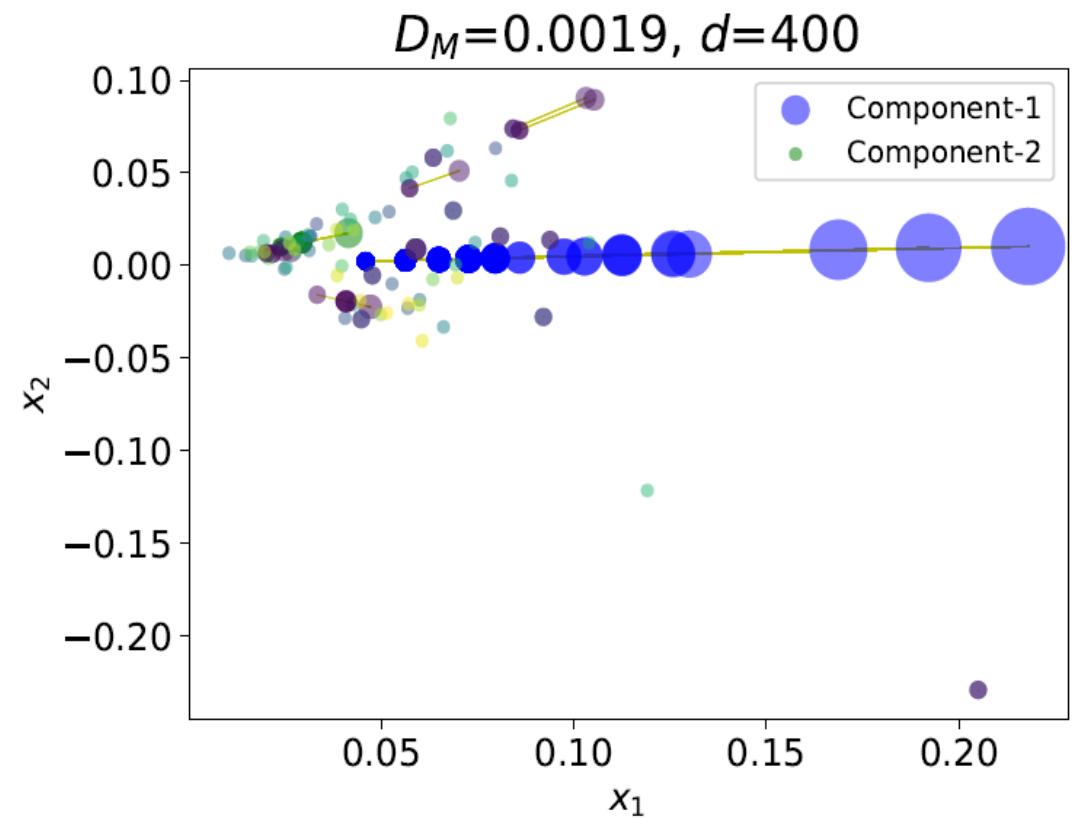
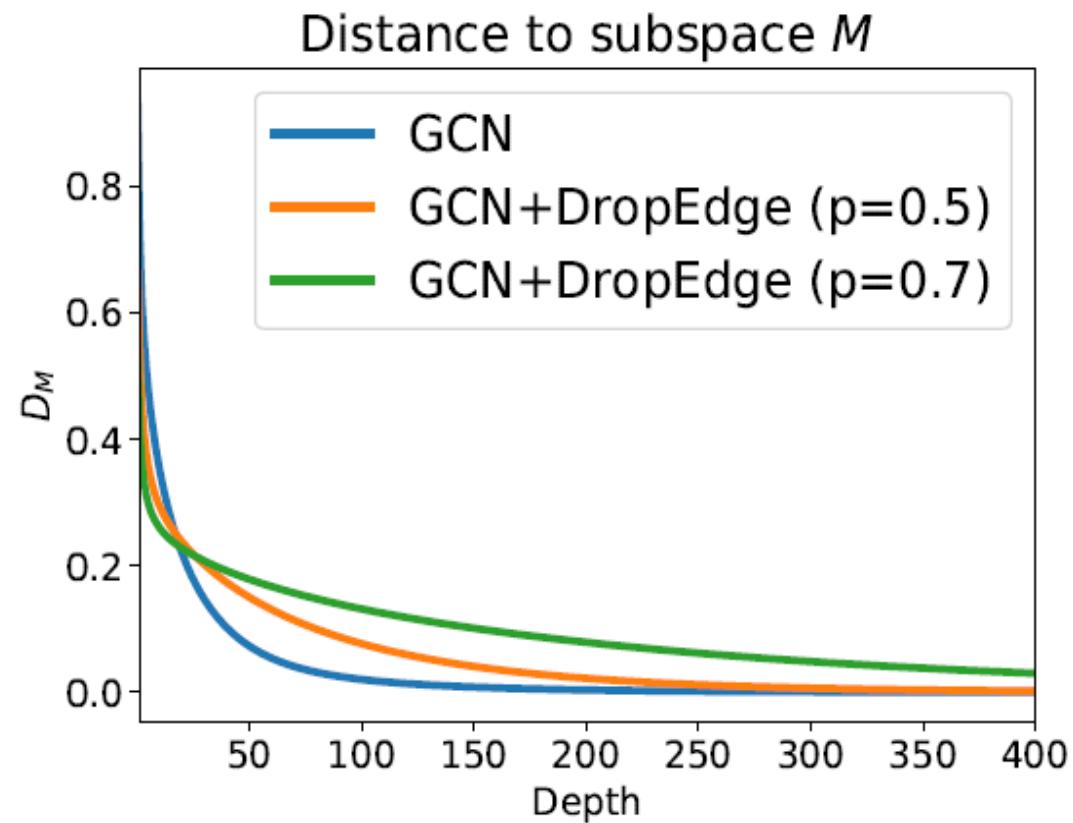
# Alleviate Over-Smoothing by Adjacency Matrix

## DropEdge results

Citeseer	4 layers	DropEdges	16 layers	DropEdges	64 layers	DropEdges
GCN	76.7	79.2(+2.5)	65.2	76.8(+11.6)	44.6	45.6(+1.0)
ResGCN	78.9	78.8(-0.1)	78.2	79.4(+1.2)	21.2	75.3(+54.1)
JKNet	79.1	80.2(+1.1)	78.8	80.1(+1.3)	76.7	80.0(+3.3)
IncepGCN	79.5	79.9(+0.4)	78.5	80.2(+1.7)	79.0	79.9(+0.9)
GraphSAGE	77.3	79.2(+1.9)	72.9	74.5(+1.6)	16.9	25.1(+8.2)
APPNP	80.3	80.8(+0.5)	80.2	81.1(+0.9)	80.4	81.3(+0.9)

# Alleviate Over-Smoothing by Adjacency Matrix

## DropEdge results



# Alleviate Over-Smoothing by Weights

Convergence speed

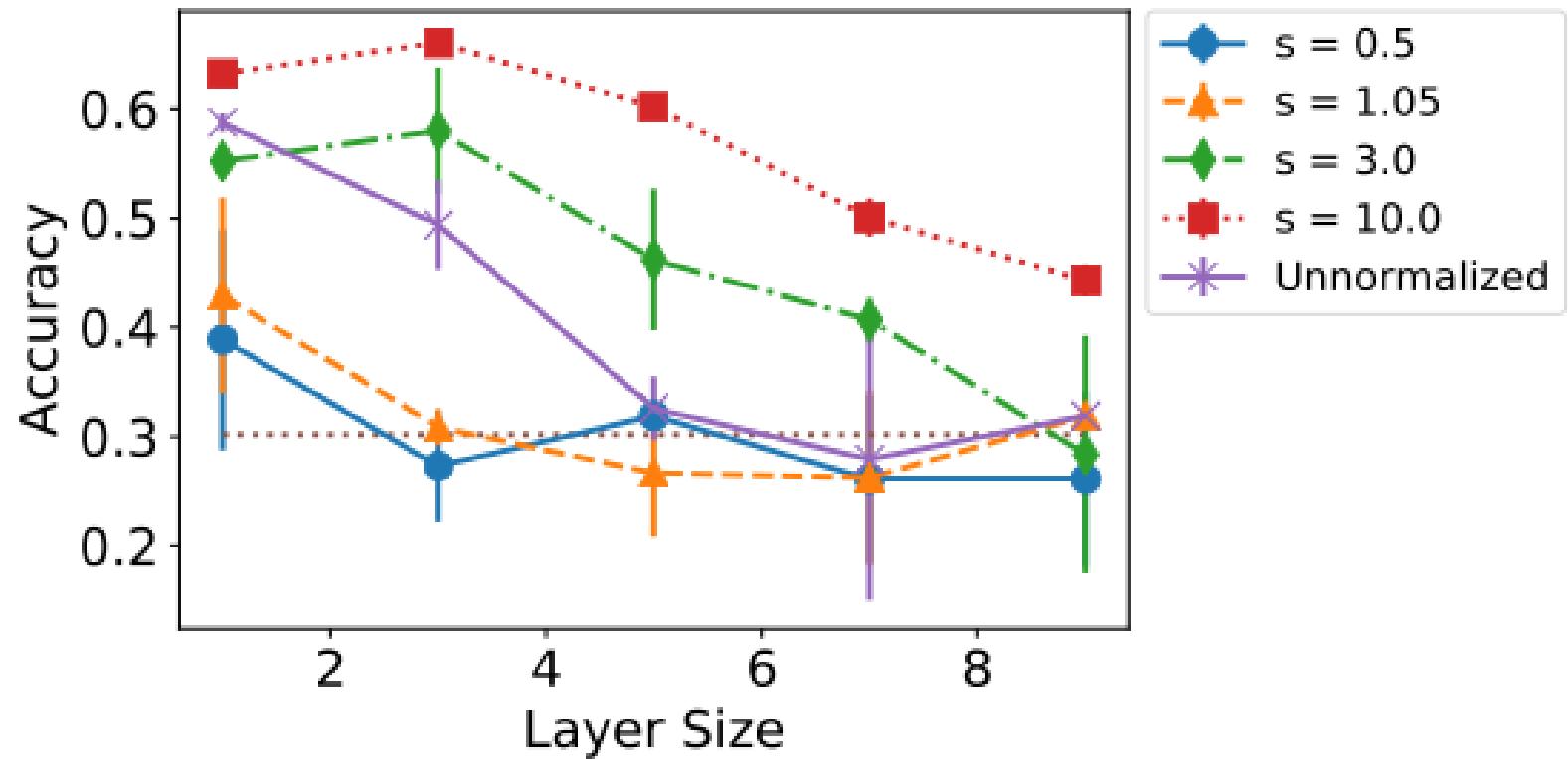
$$\nu = \lambda_{m+1} s$$

Weights

- Similarly, increasing  $s$  will also increase  $\nu$ . So how to increase  $s$ ? Increase the initial  $W_l$ s.

# Alleviate Over-Smoothing by Weights

Try different  $s$  as initial



# Training deep GNNs

- ◀ Why do we need deep GNNs?
- ◀ Can GNNs simply go deep?
- ◀ What impedes GNNs to go deep?
  - ◀ Overfitting (Common)
  - ◀ Training dynamics (Common)
  - ◀ Over-smoothing (Graph Specific)
- ◀ How to make GNNs deep?
  - ◀ Architecture refinement
  - ◀ Manipulating input (DropEdge)
  - ◀ Layer normalizations

# Pair Norm: Center and Rescale

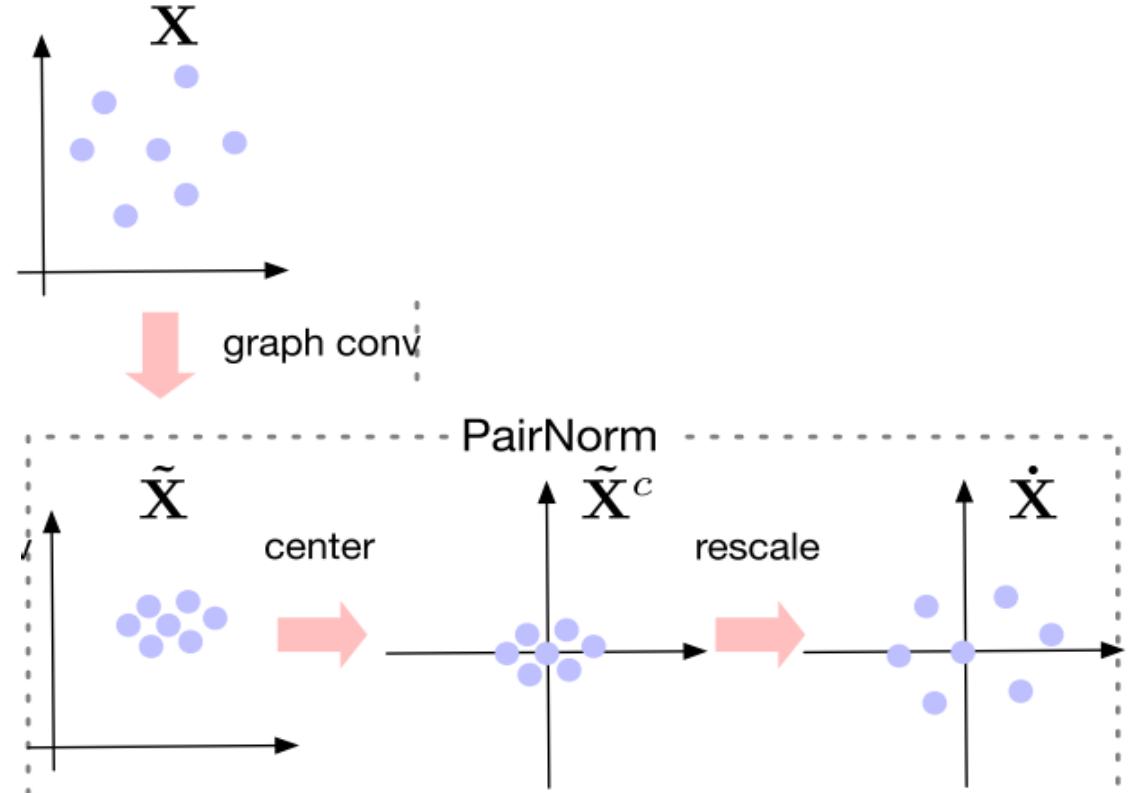
PairNorm: Center and rescale (normalize) GCN outputs  $\tilde{X} := \text{GCN}(A, X)$  to keep the **total pairwise squared distance unchanged**

Center

$$\tilde{x}_i^c = \tilde{x}_i - \frac{1}{n} \sum_{i=1}^n \tilde{x}_i$$

Rescale

$$\dot{x}_i = s\sqrt{n} \frac{\tilde{x}_i^c}{\sqrt{\|\tilde{X}^c\|_F^2}}$$



See also GroupNorm (Zhou et al., 2020)



Tencent  
AI Lab



# Break



5 min



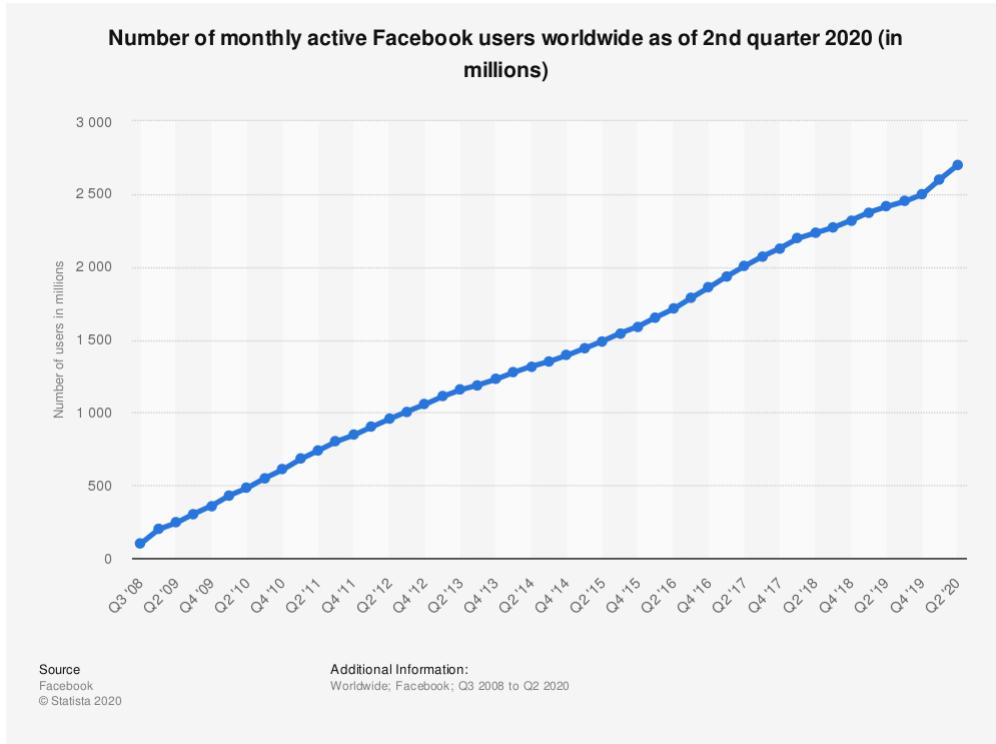
Tencent  
AI Lab



# Scalability of GNNs

# Graph in the Real-world Can be Very Large

- Large scale:



Large number:

ZINC Substances Catalogs Tranches Biological More

## ZINC15

Welcome to ZINC, a free database of commercially-available compounds for virtual screening. ZINC contains over 230 million purchasable compounds in ready-to-dock, 3D formats. ZINC also contains over 750 million purchasable compounds you can search for analogs in under a minute.

### Getting Started

- Getting Started
- What's New
- About ZINC 15 Resources
- Current Status / In Progress
- Why are ZINC results "estimates"?

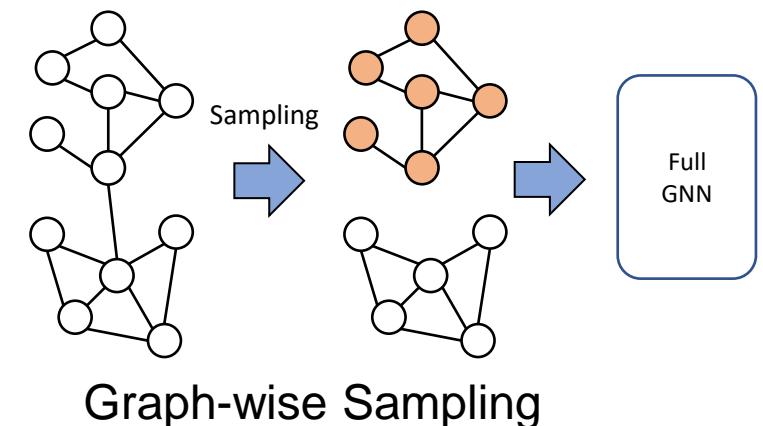
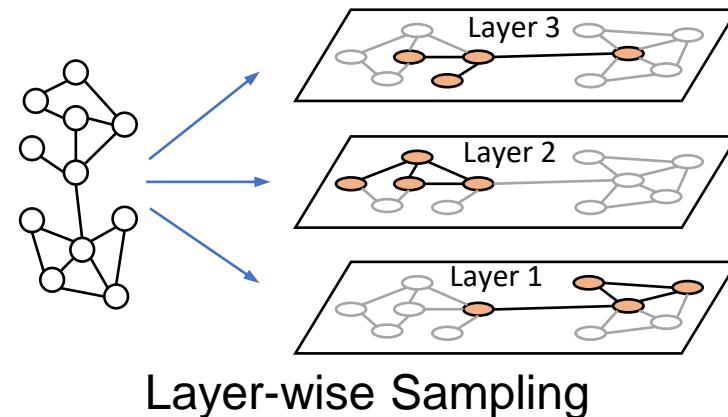
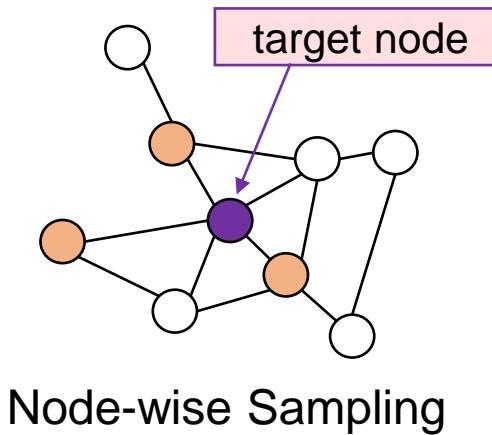
### Ask Questions

You can use ZINC for **general** questions such as

- How many substances in current clinical trials have PAINS
- How many natural products have names in ZINC and are SMILES, names and calculated logP
- How many endogenous human metabolites are there? (

# Three Paradigms

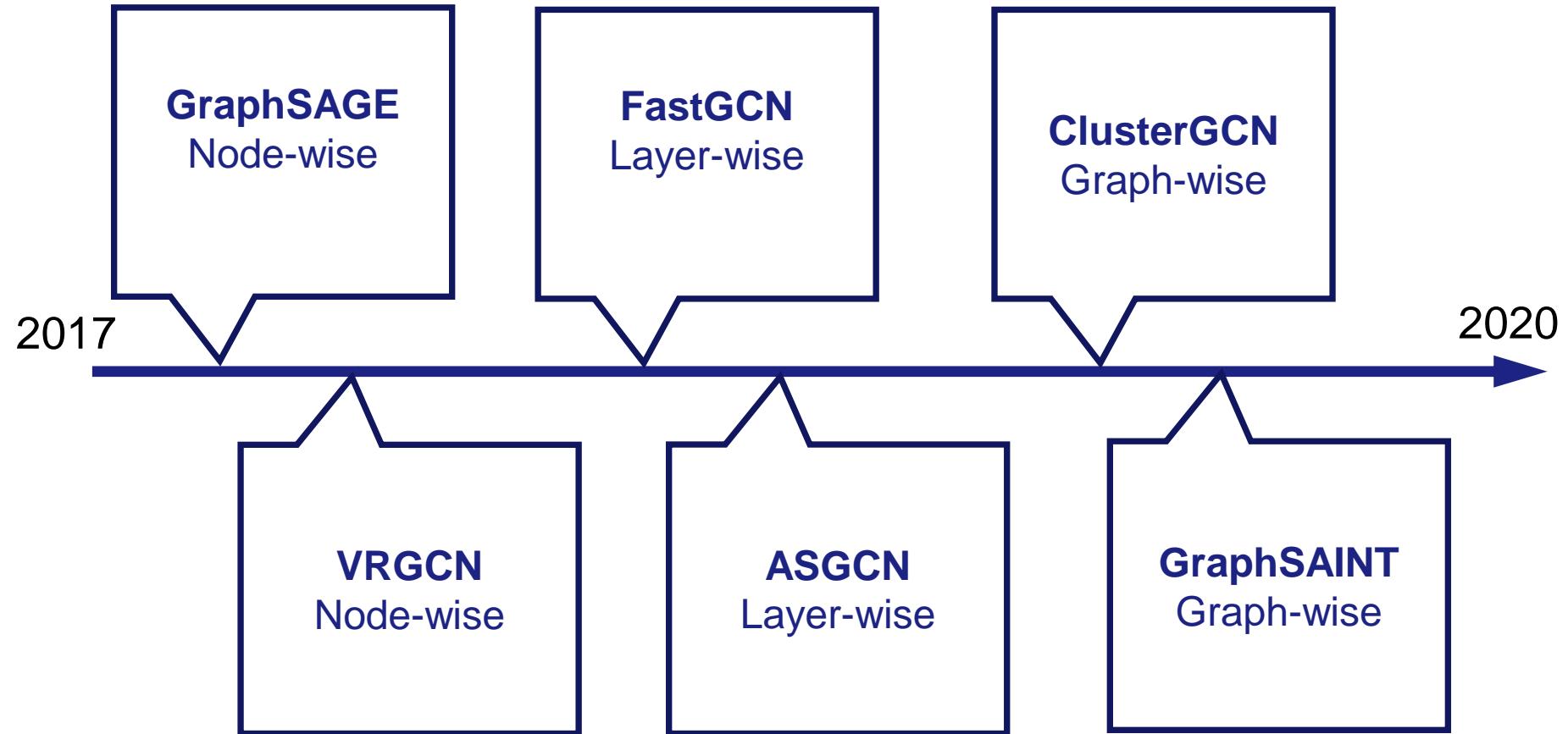
- Why the original GNN fails on large graph?
  - Large memory requirement.
  - Inefficient gradient update.
- Three paradigms toward large-scale GNN:



# Two issues towards the large-scale GNNs

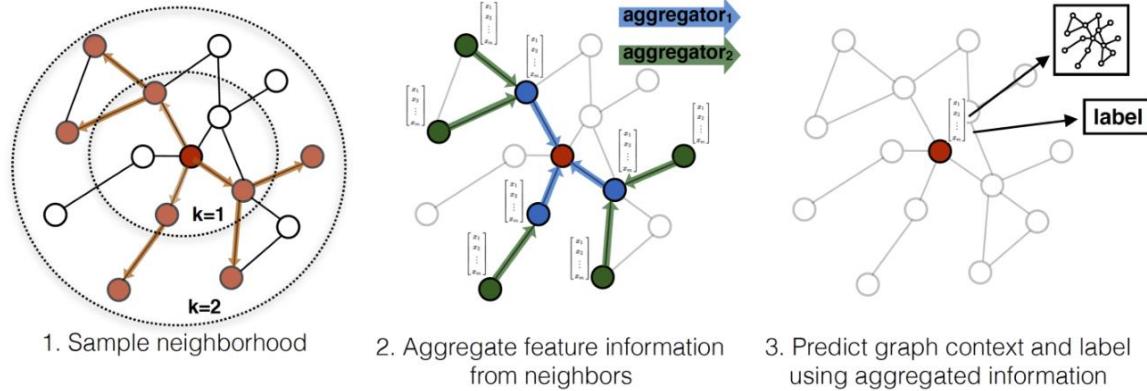
- How to design efficient sampling algorithm?
- How to guarantee the sampling quality?

# Overview





# Node-wise Sampling: GraphSAGE



- The architecture:

$$\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$$

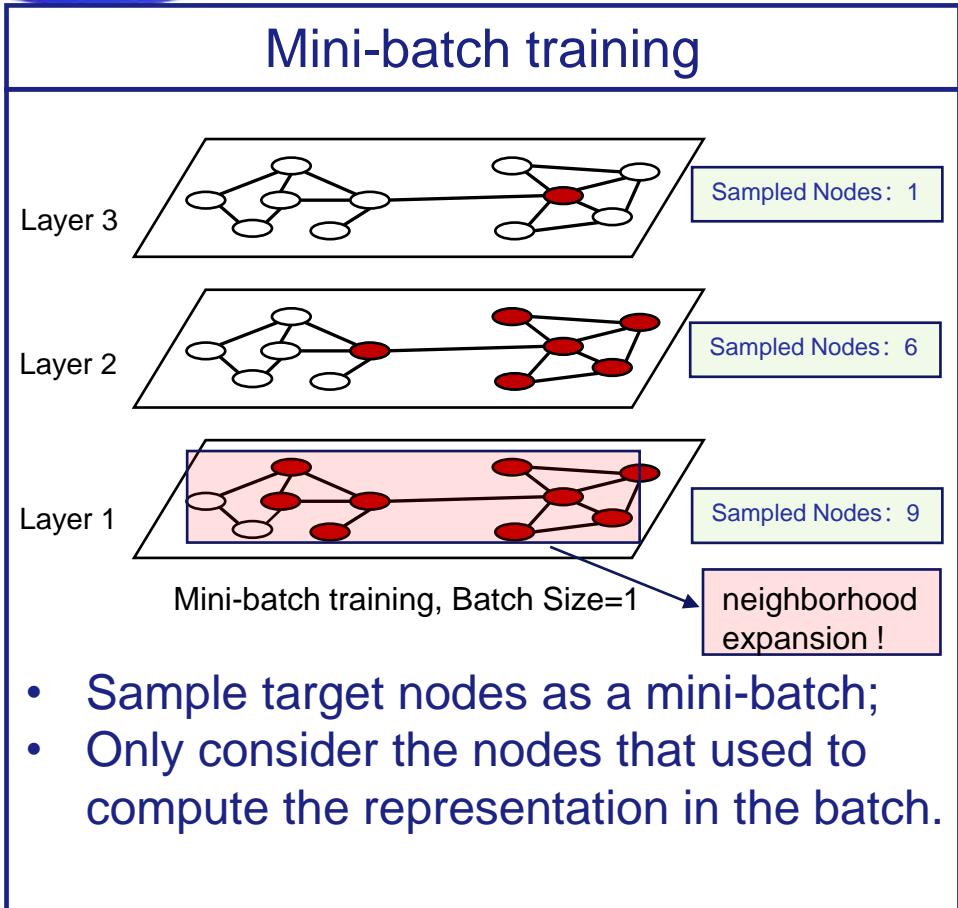
$$\mathbf{h}_v^k \leftarrow \sigma(\mathbf{W}^k \cdot \text{CONCAT}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}))$$

Use Concatenation instead of SUM

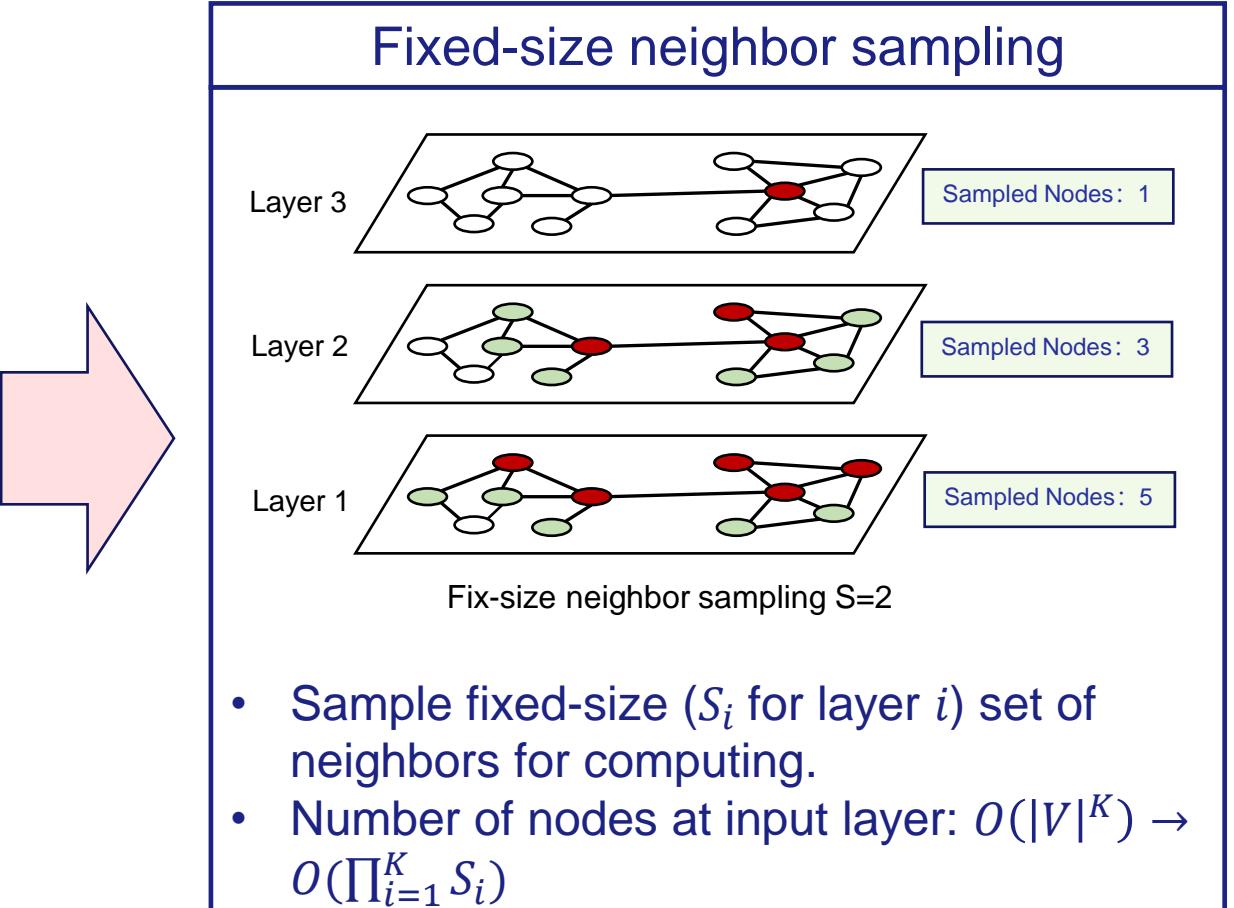
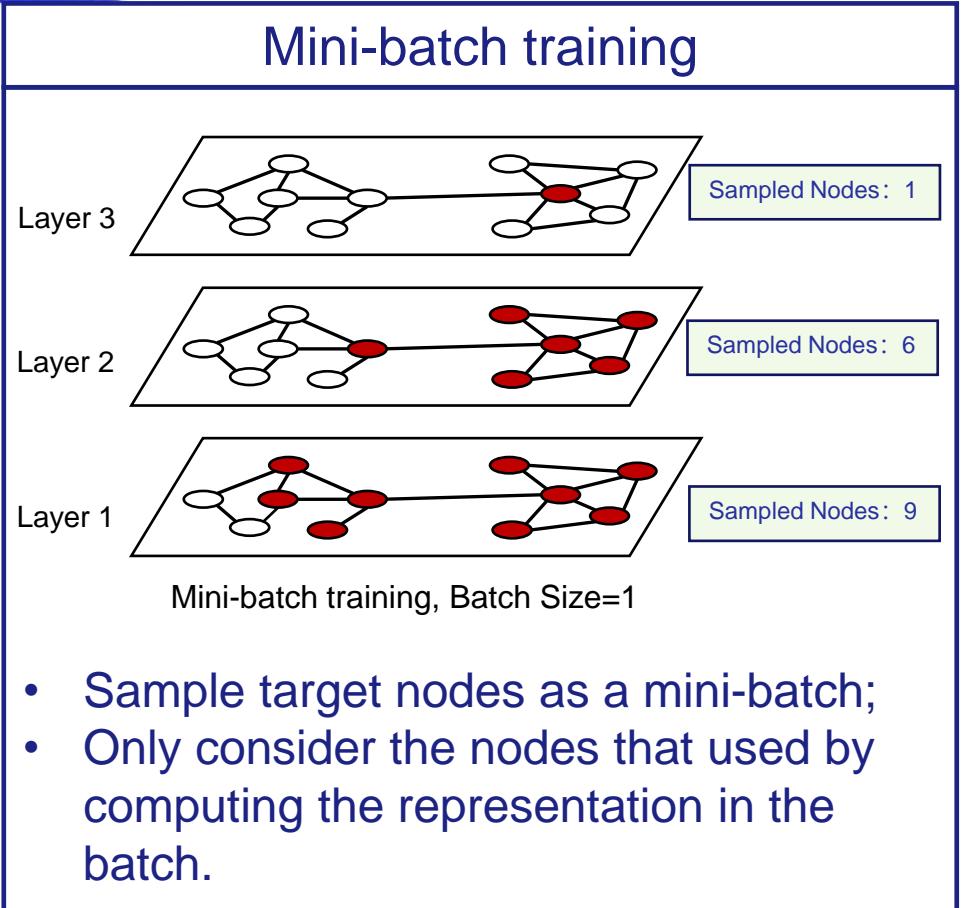
Generalized Aggregator

- Mean aggregator (GCN)
- Pooling aggregator
- LSTM aggregator
- .....

# Towards large-scale GraphSAGE



# Towards large-scale GraphSAGE

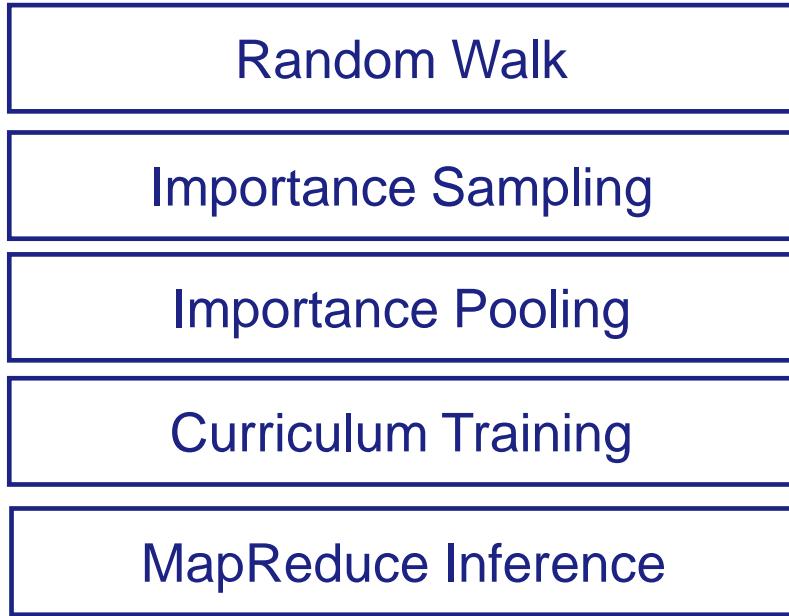
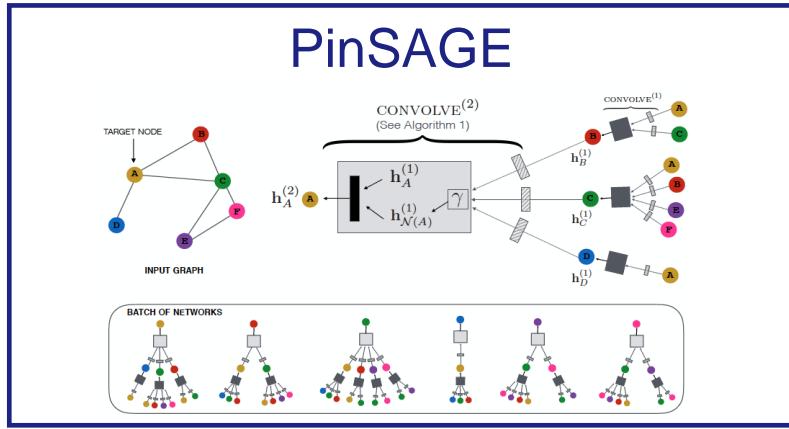
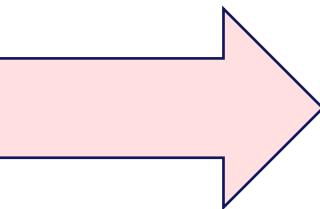
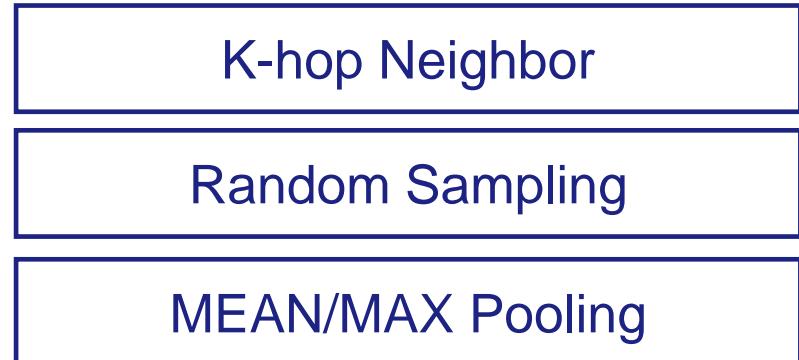
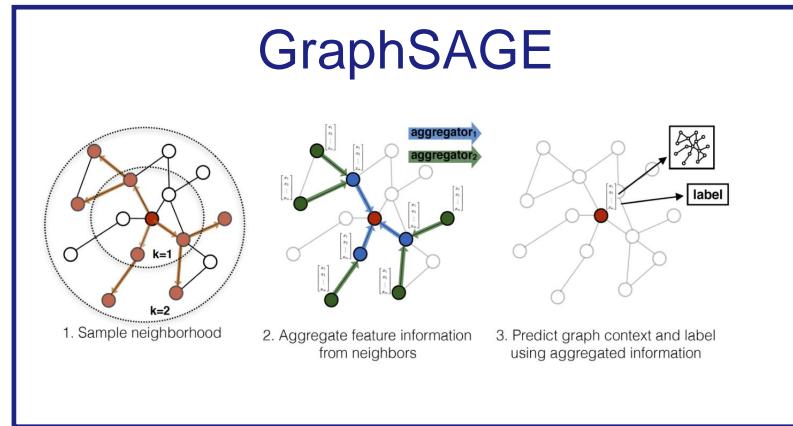


# Node-wise Sampling: GraphSAGE

- Pros:
  - Generalized aggregator.
  - Mini-batch training and fixed-size neighbor sampling.
- Cons:
  - Neighborhood expansion on deeper GNNs.
  - No guarantees for the sampling quality.

# GraphSAGE Extension: PinSAGE

- An early-stage industrial GNN-based recommendation system.



# Node-wise Sampling: VR-GCN

- GraphSAGE: Neighborhood Sampling

$$\text{NS}_u^{(l)} := R \sum_{v \in \hat{\mathcal{N}}^{(l)}(u)} P_{uv} \mathbf{h}_v^{(l)}, \quad R = |\mathcal{N}(u)|/D^{(l)}$$

- biased sampling / larger variance  $\rightarrow$  larger sample size  $\hat{\mathcal{N}}^{(l)}(u)$ .

- Control Variate Based Estimator (**CV Sampler**):

- Maintain the historical hidden embedding  $\bar{\mathbf{h}}_v^{(l)}$  for a better estimation.
- **Variance reduction  $\rightarrow$  Variance elimination  $\rightarrow$  Smaller sample size  $\hat{\mathcal{N}}^{(l)}(u)$ .**

- VR-GCN:

$$\mathbf{H}^{(l+1)} = \sigma(\widehat{\mathbf{P}}^{(l)}(\mathbf{H}^{(l+1)} - \bar{\mathbf{H}}^{(l)}) + \mathbf{P}\bar{\mathbf{H}}^{(l)})\mathbf{W}^{(l)}$$

The sampled normalized adjacency matrix at layer  $l$ .

Historical hidden embedding

- **One more thing:** CVD Sampler -- Control Variate for Dropout.

$\mathcal{N}(u)$ : the neighbor set of node  $u$ .

$\hat{\mathcal{N}}^{(l)}(u) \subset \mathcal{N}(u)$ : the sampled neighbor set of node  $u$  at layer  $l$ .

$\mathbf{P}$ : the normalized adjacency matrix.

$\widehat{\mathbf{P}}^{(l)}$ : the sampled normalized adjacency matrix at layer  $l$ .

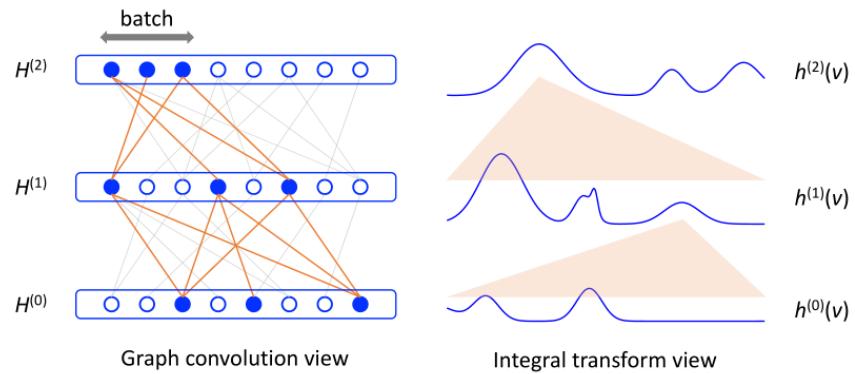
$D^{(l)}$ : the sampled size for each node at layer  $l$



# Node-wise Sampling: VR-GCN

- Pros:
  - Analyze the variance reduction on node-wise sampling.
  - Successfully reducing the size of samples.
- Cons:
  - Additional memory consuming for storing the historical hidden embeddings.

# Layer-wise Sampling: FastGCN



The functional generalization of GCN

$$L = E_{v \sim P} [g(h^{(M)}(v))] = \int g(h^{(M)}(v)) dP(v)$$

Integral Transform



The estimation

$$L_{t_0, t_1, \dots, t_M} := \frac{1}{t_M} \sum_{i=1}^{t_M} g(h_{t_M}^{(M)}(u_i^{(M)}))$$

The i.i.d. node sample at each layer  
 → bootstrapping

**FastGCN:** sampling fixed number of nodes at each layer.

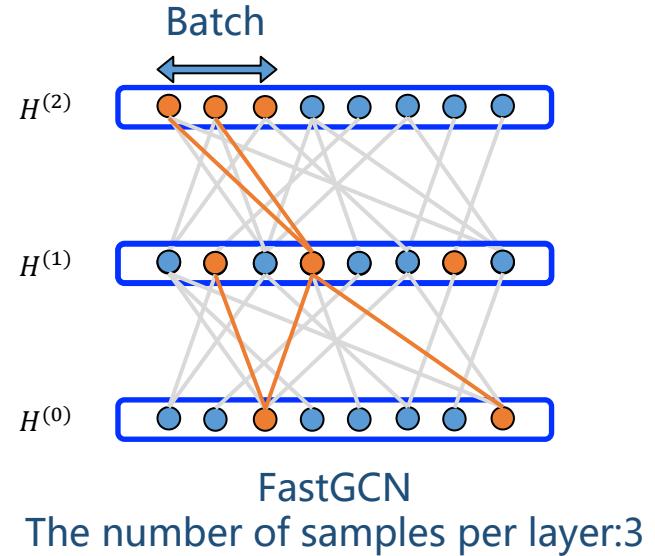
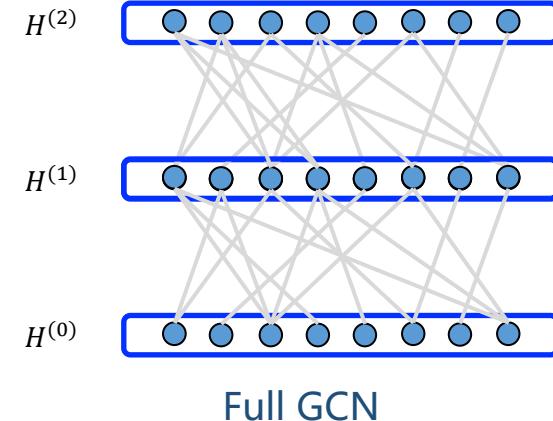
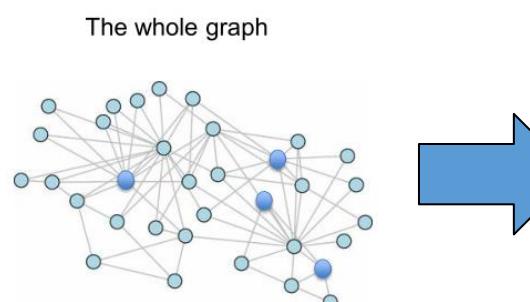


# Layer-wise Sampling: FastGCN

- Towards the variance reduction: importance sampling.

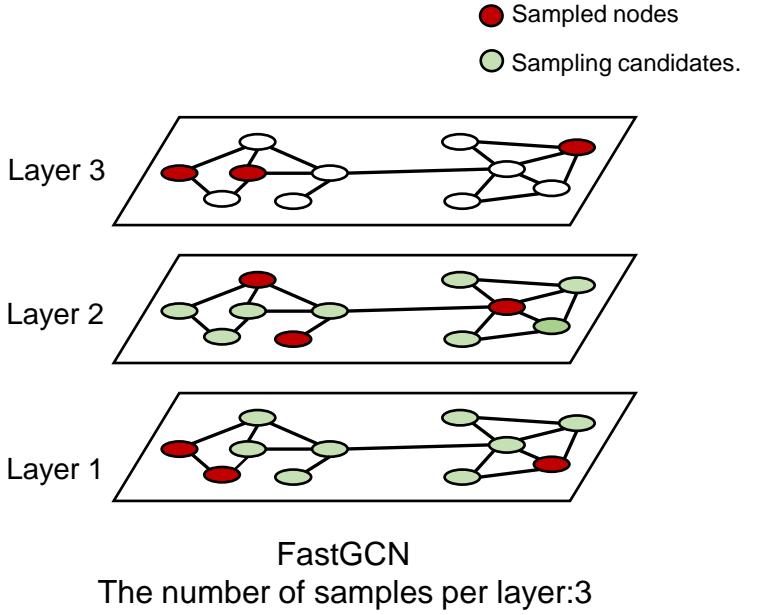
$$q(u) = \|\hat{A}(:, u)\|^2 / \sum_{u' \in V} \|\hat{A}(:, u')\|^2, \quad u \in V$$

This sampling probability keeps the same for each layer.



# Layer-wise Sampling: FastGCN

- Pros:
  - Avoid neighborhood expansion problem.
  - Sample method with quality guarantee.
- Cons:
  - Failed to capture the between-layer correlations.
  - Performance compromise.





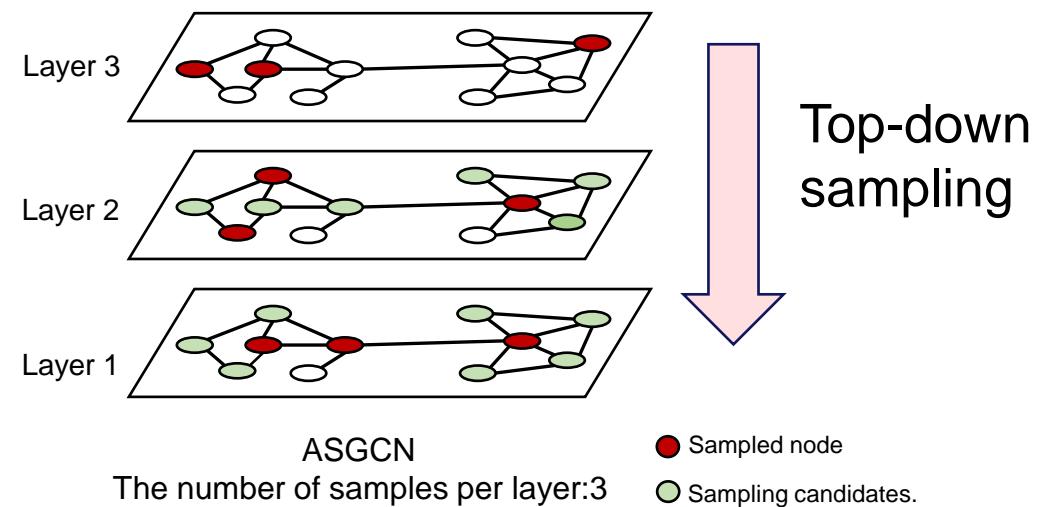
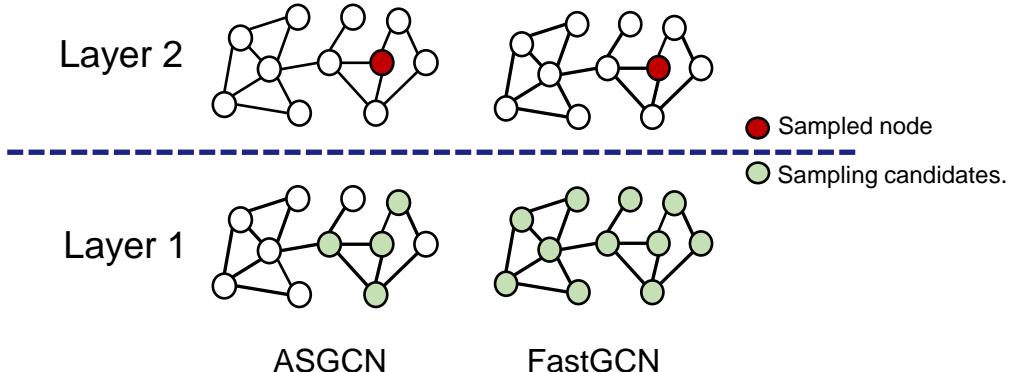
# Layer-wise Sampling: ASGCN

- Adaptive layer-wise sampling:
  - The sampling probability of lower layers depends on the upper ones.

$$q(u_j) = \frac{p(u_j|v_i)}{p(u_j|v_1 \dots v_n)} \quad s.t. \quad p(u_j|v_i) = \frac{\hat{a}(v_i, u_j)}{N(v_i)}, N(v_i) = \sum_{j=1}^n \hat{a}(v_i, u_j)$$

the probability of sampling node  $u_j$  given node  $v_i$ .

The entry of node  $v_i$  and  $u_j$  in re-normalization of the adjacency matrix  $\hat{A}$ .



$x(u_j)$ : the node feature of node  $u_j$ .

$p(u_j|v_i)$ : the probability of sampling node  $u_j$  given node  $v_i$ .

$\hat{a}(v_i, u_j)$ : The entry of node  $v_i$  and  $u_j$  in re-normalization of the adjacency matrix  $\hat{A}$ .

$\hat{u}_q(v_i)$ : the output hidden embeddings of node  $v_i$ .

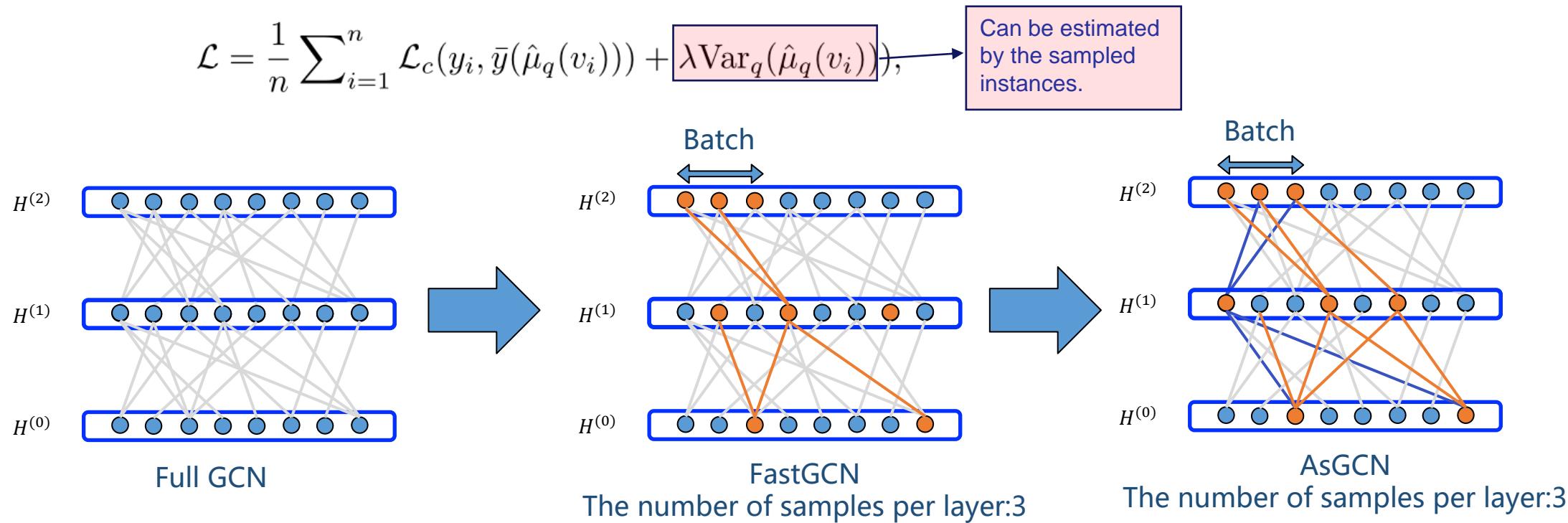
# Layer-wise Sampling: ASGCN

- Parameterized for **explicit variance reduction.**

$$q^*(u_j) = \frac{\sum_{i=1}^n p(u_j|v_i)|g(x(u_j))|}{\sum_{j=1}^N \sum_{i=1}^n p(u_j|v_i)|g(x(v_j))|}, \quad g(x(u_j)) = W_g x(u_j)$$

$x(u_j)$ : the node feature of node  $u_j$ .  
 $p(u_j|v_i)$ : the probability of sampling node  $u_j$  given node  $v_i$ .  
 $\hat{a}(v_i, u_j)$ : The entry of node  $v_i$  and  $u_j$  in re-normalization of the adjacency matrix  $\hat{A}$ .  
 $\hat{\mu}_q(v_i)$ : the output hidden embeddings of node  $v_i$ .

- Optimize the sampler  $q^*(u_j)$  to minimize the variance:



# Layer-wise Sampling: ASGCN

- Self-dependent shared attention:

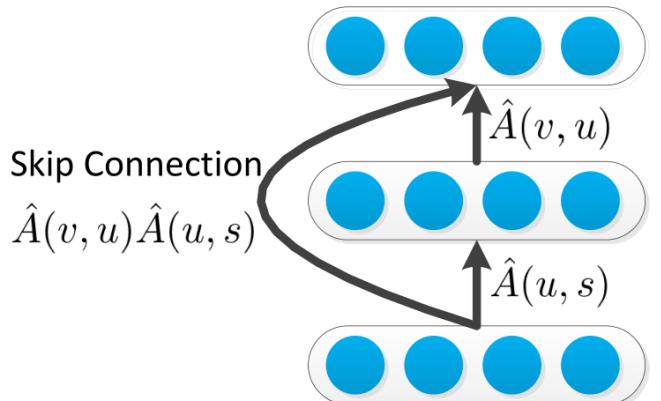
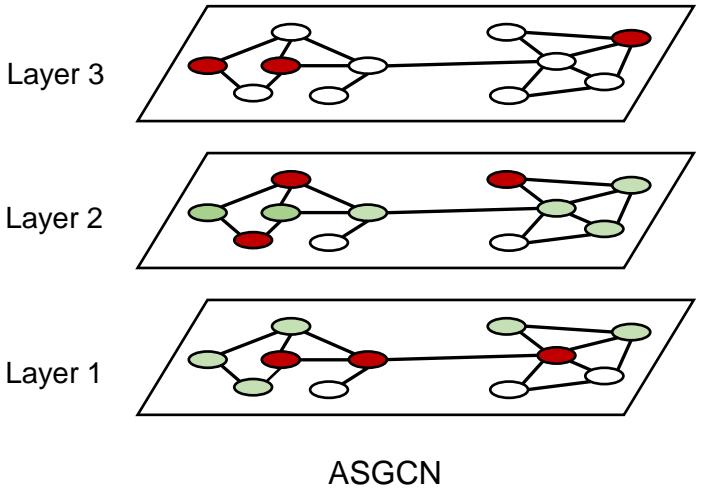
$$a(x(v_i), x(u_j)) = \frac{1}{n} \text{ReLU}(W_1 g(x(v_i)) + W_2 g(x(u_j)))$$

- Preserving second-order proximities by skip connections:

$$h_{skip}^{(l+1)}(v_i) = \sum_{j=1}^n \hat{a}_{skip}(v_i, s_j) h^{(l-1)}(s_j) W_{skip}^{(l-1)}, \quad i = 1, \dots, n,$$

- where,  $\hat{a}_{skip}$  is estimated by:

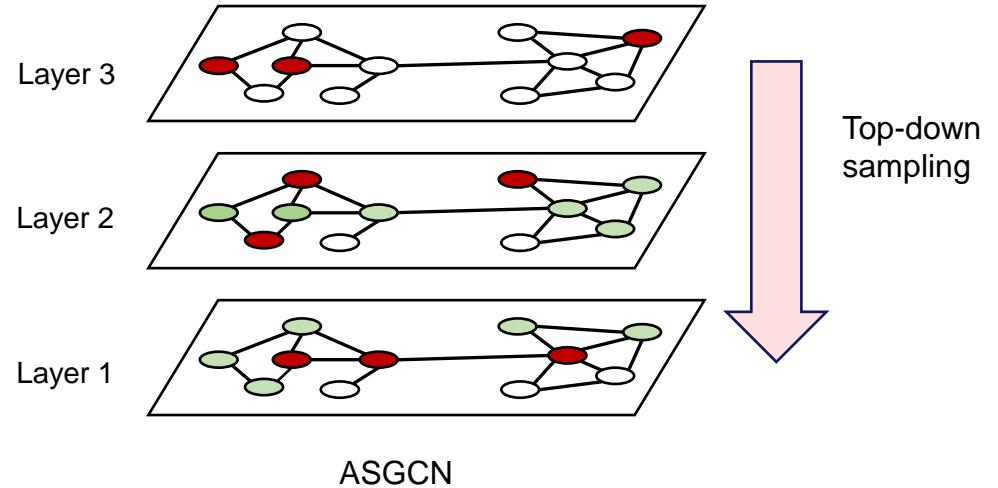
$$\hat{a}_{skip}(v_i, s_j) \approx \sum_{k=1}^n \hat{a}(v_i, u_k) \hat{a}(u_k, s_j).$$





# Layer-wise Sampling: ASGCN

- Pros:
  - Good performance.
  - Better variance control.
- Cons:
  - Additional dependence during sampling.



ASGCN

Github: <https://github.com/huangwb/AS-GCN>

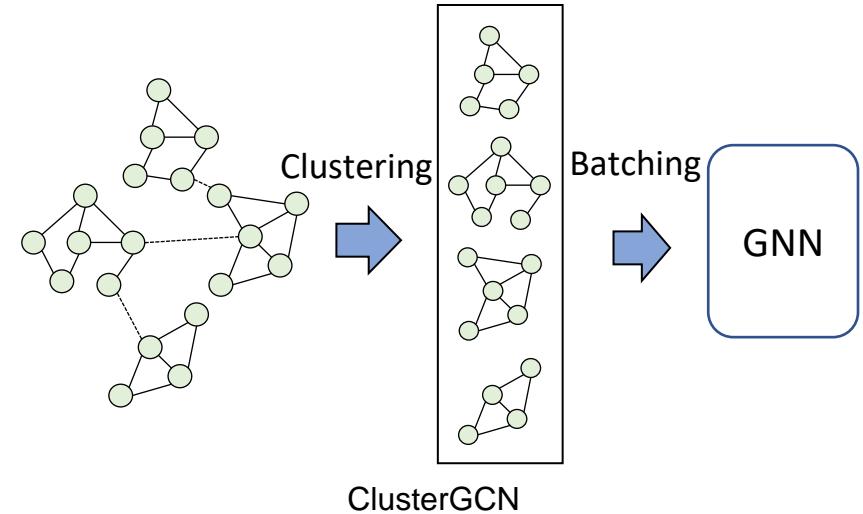
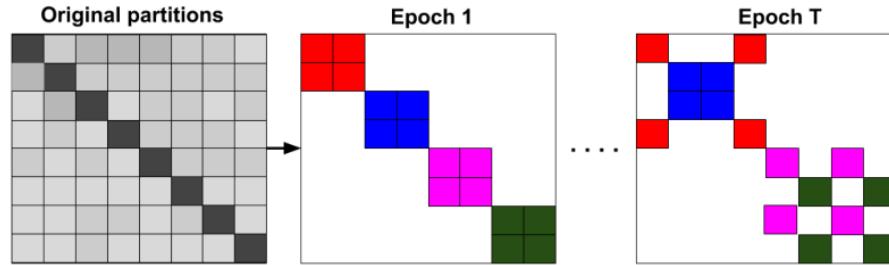
# Graph-wise Sampling: ClusterGCN

- Extract small clusters based efficient clustering algorithms.

$$\bar{G} = [G_1, \dots, G_c] = [\{\mathcal{V}_1, \mathcal{E}_1\}, \dots, \{\mathcal{V}_c, \mathcal{E}_c\}],$$

$$\bar{A} = \begin{bmatrix} A_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & A_{cc} \end{bmatrix}, \Delta = \begin{bmatrix} 0 & \cdots & A_{1c} \\ \vdots & \ddots & \vdots \\ A_{c1} & \cdots & 0 \end{bmatrix},$$

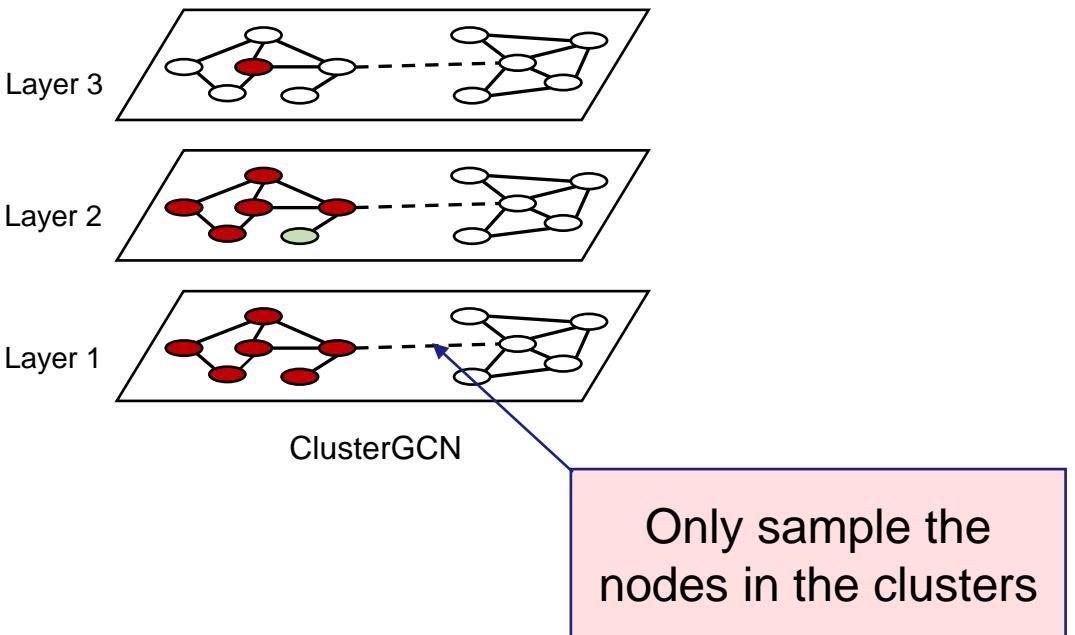
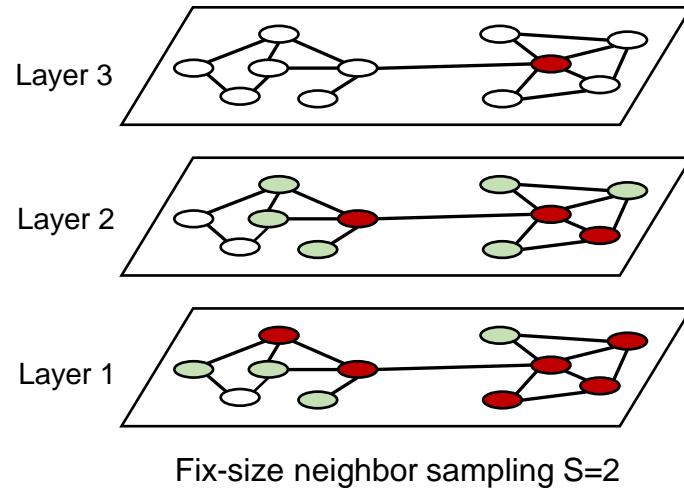
- Random batching at the subgraph level.





# Graph-wise Sampling: ClusterGCN

- Neighbor expansion control.



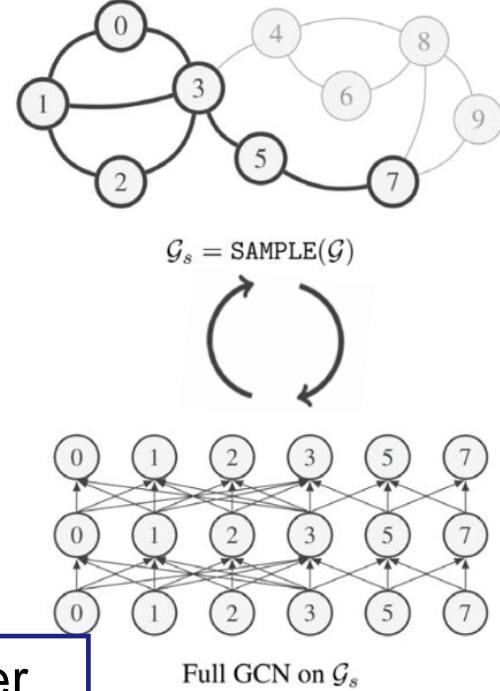
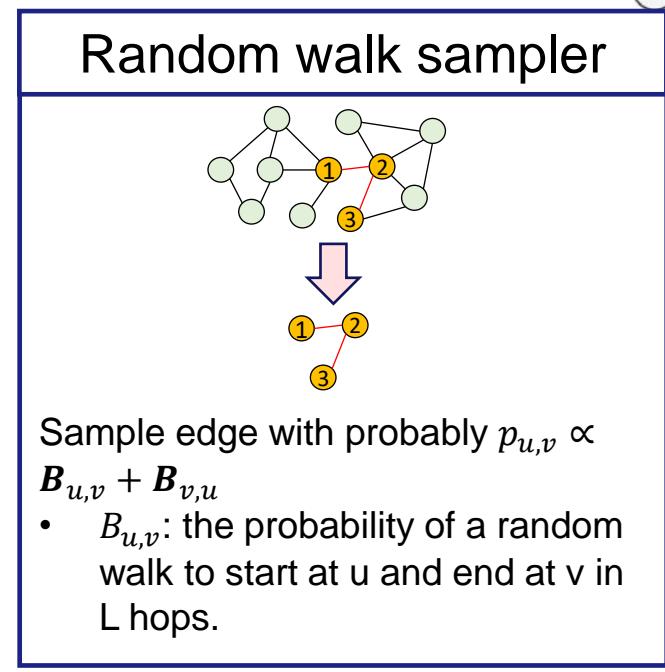
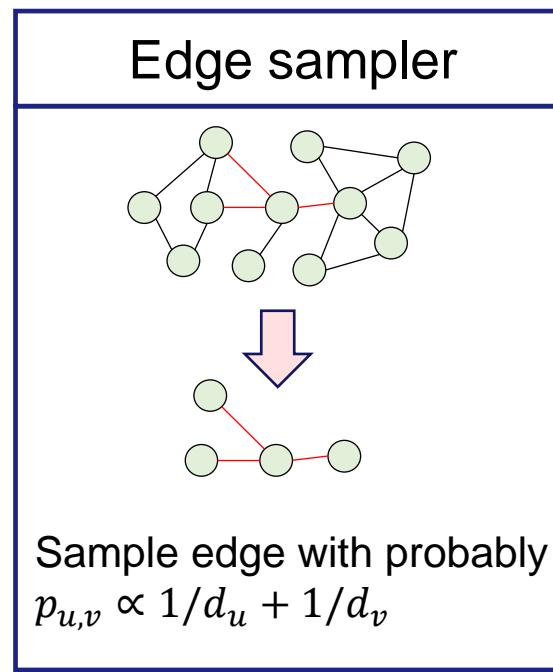
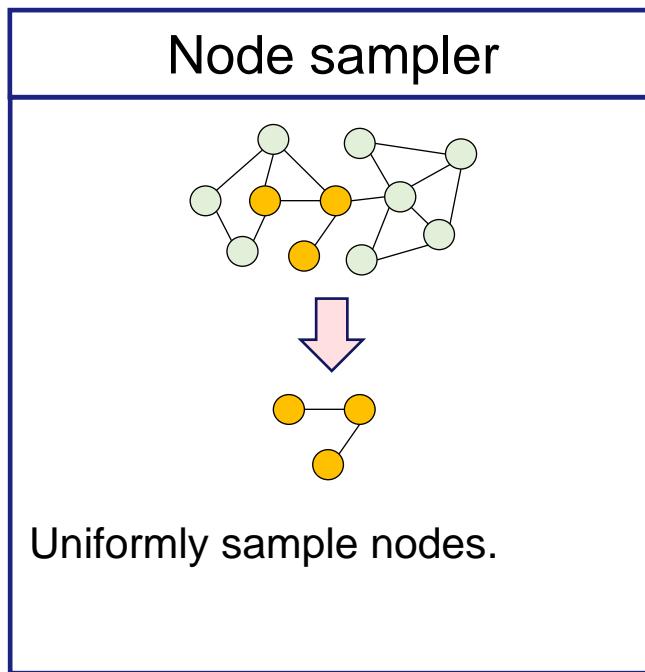


# Graph-wise Sampling: ClusterGCN

- Pros:
  - Good performance / Good memory usage.
  - Alleviate the neighborhood expansion problem in traditional mini-batch training.
- Cons:
  - Empirical results without analyzing the sampling quality.

# Graph-wise Sampling: GraphSAINT

- Directly sample a subgraph for mini-batch training according to subgraph sampler.
- Sampler construction



# Graph-wise Sampling: GraphSAINT

- How to eliminate the bias introduce by the sampler?

- Loss normalization:

$$\mathcal{L}_{\text{batch}} = \sum_{v \in G_s} L_v / \lambda_v, \lambda_v = |V| p_v.$$

- Aggregation normalization:

$$a(u, v) = p_{u,v} / p_v$$

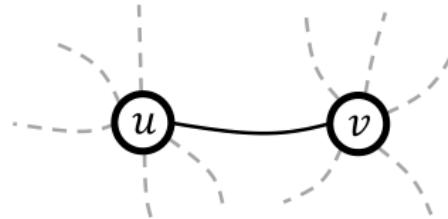
$p_v$ : the probability of a node  $v \in V$  being sampled.

$p_{u,v}$ : the probability of an edge  $(u, v) \in E$  being sampled.

# Graph-wise Sampling: GraphSAINT

- How to reduce the sampling variance?
  - The optimal edge probability for variance minimization:

$$p_{u,v} \propto 1/d_u + 1/d_v$$



Sample "—" more  
Sample "—" less

$p_v$ : the probability of a node  $v \in V$  being sampled.  
 $p_{u,v}$ : the probability of a edge  $(u, v) \in E$  being sampled.

# Graph-wise Sampling: GraphSAINT

Dataset	Nodes	Edges	Degree	Feature	Classes	Train / Val / Test
PPI	14,755	225,270	15	50	121 (m)	0.66 / 0.12 / 0.22
Flickr	89,250	899,756	10	500	7 (s)	0.50 / 0.25 / 0.25
Reddit	232,965	11,606,919	50	602	41 (s)	0.66 / 0.10 / 0.24
Yelp	716,847	6,977,410	10	300	100 (m)	0.75 / 0.10 / 0.15
Amazon	1,598,960	132,169,734	83	200	107 (m)	0.85 / 0.05 / 0.10

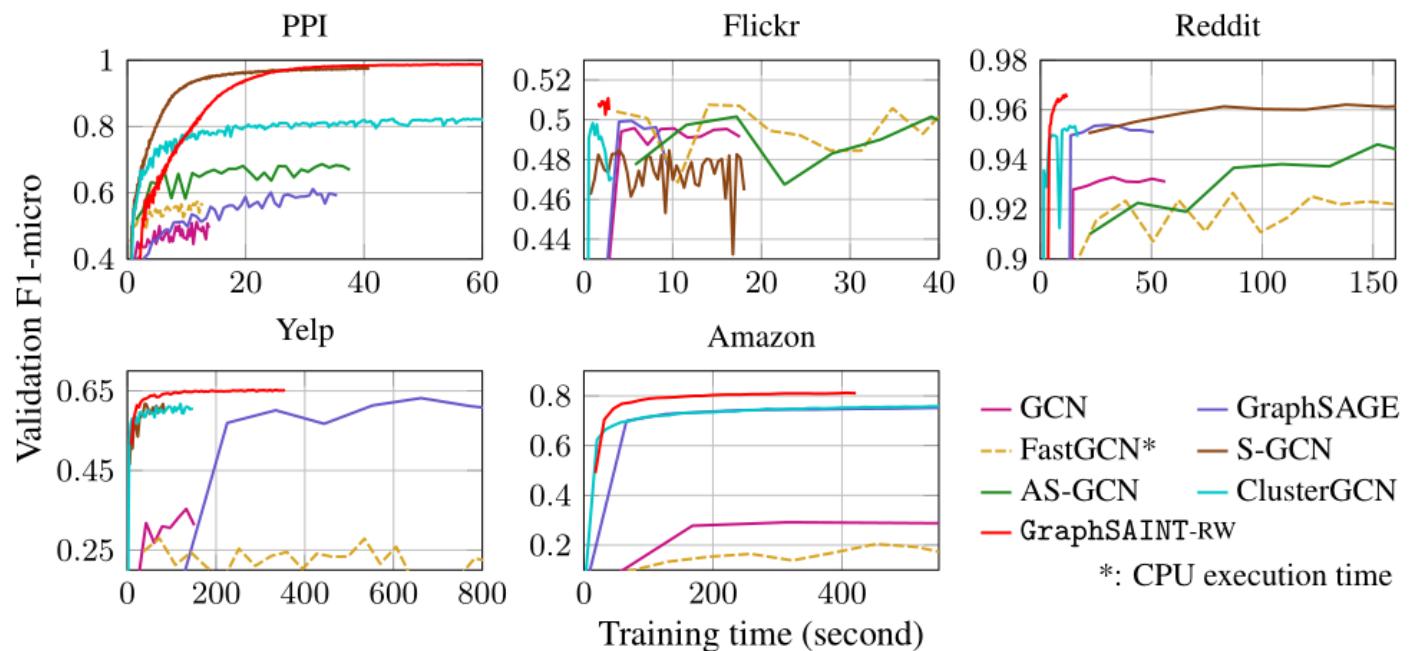


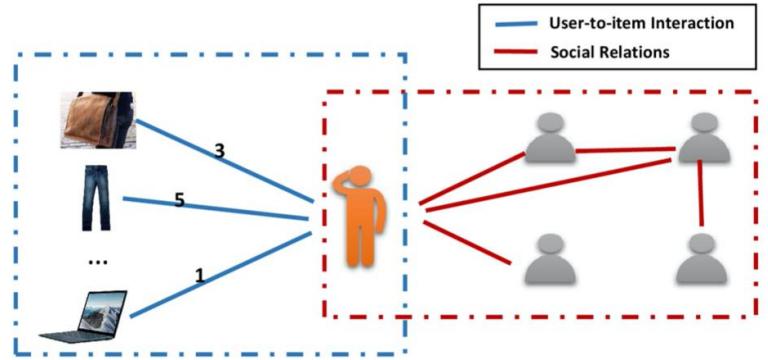
Figure 2: Convergence curves of 2-layer models on GraphSAINT and baselines

- Highly flexible and extensible (graph samplers, GNN architectures, etc. )
- Good performance (accuracy, speed)

# Large-Scale GNN in Real-World Applications

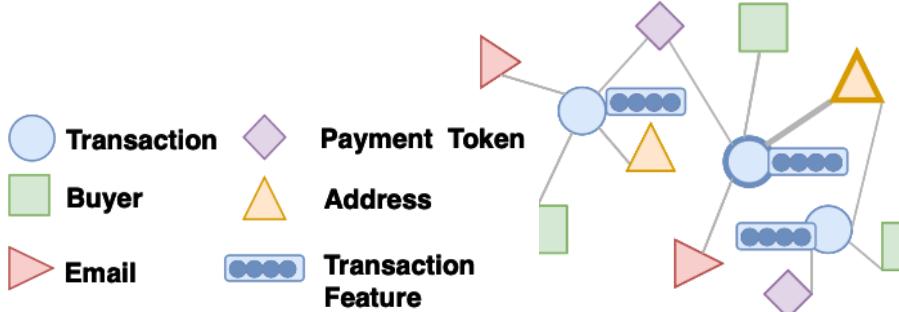
- Large-scale GNN has been successfully applied to real-world applications.

## Social Recommendation



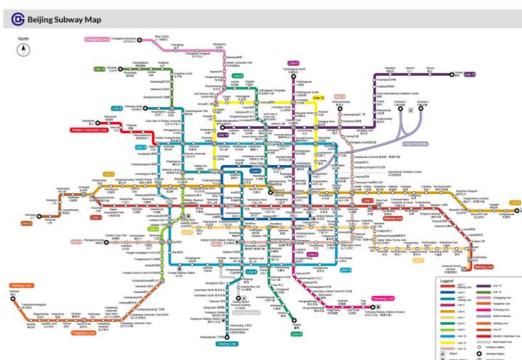
Fan W, Ma Y, Li Q, et al.

## Fraud Transaction Detection



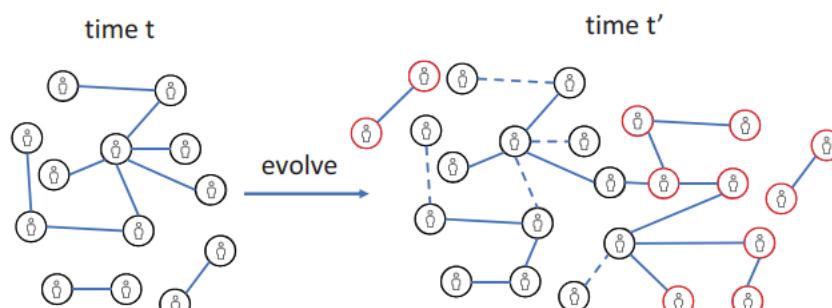
Rao S X, Zhang S, Han Z, et al. . 2020

## Traffic Forecasting



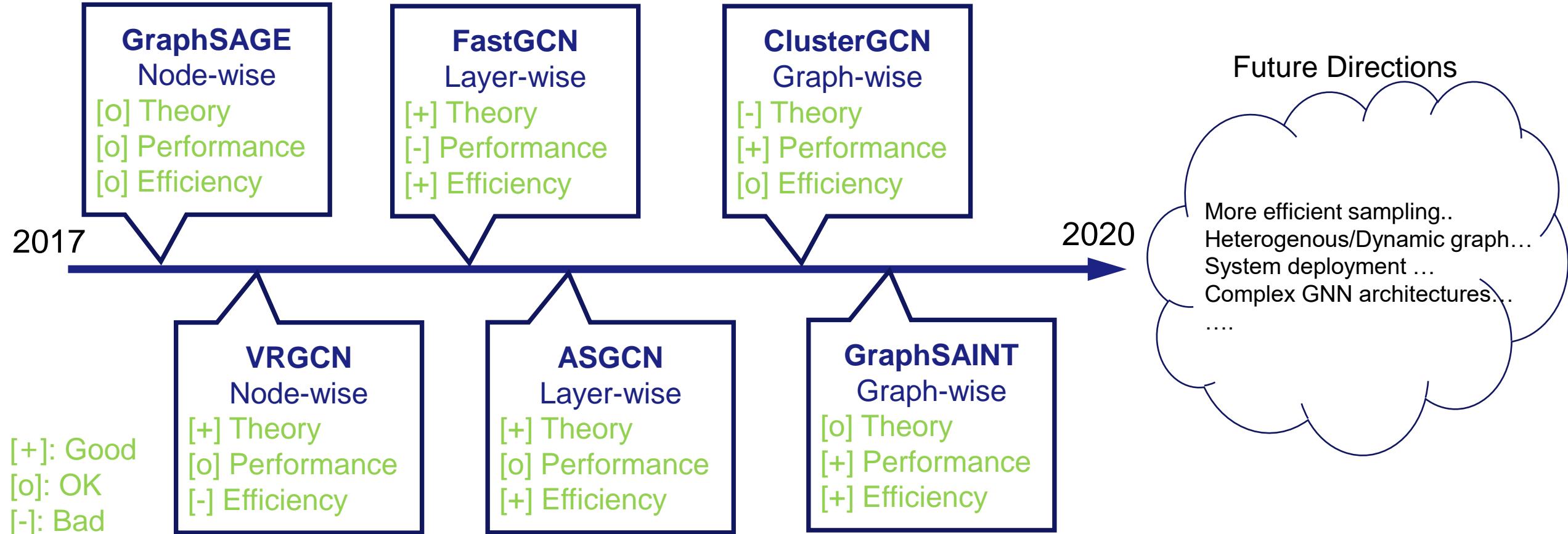
Wang X, Ma Y, Wang Y, et al. 2021

## Bitcoin Transaction Forecasting



Wei W, Zhang Q, Liu L. 2020

# Summary





Tencent  
AI Lab



# Robustness of GNNs



# Adversarial Attacks on Deep Learning



Classified as panda

$x$

Small adversarial noise

$\epsilon$

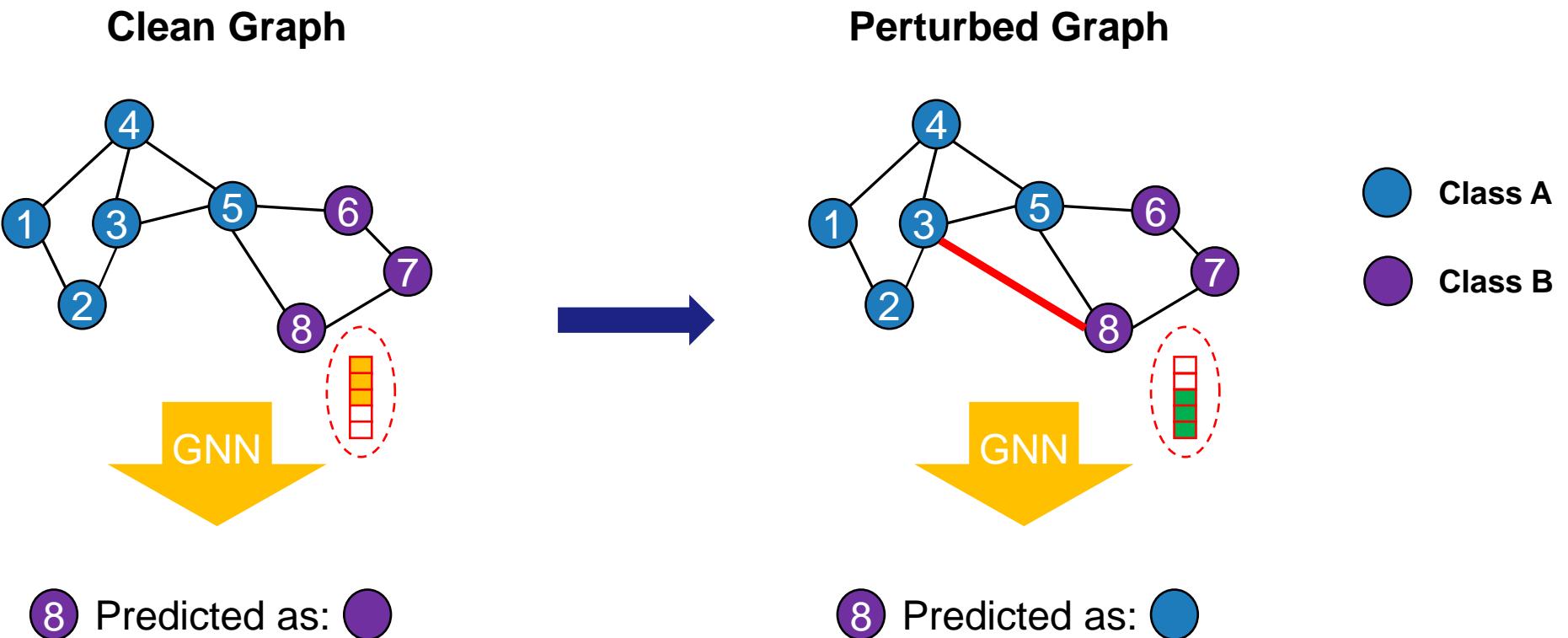
Classified as gibbon

$x'$

Find  $x'$  satisfying  $\|x' - x\| \leq \Delta$  s.t.  $C(x') \neq y$

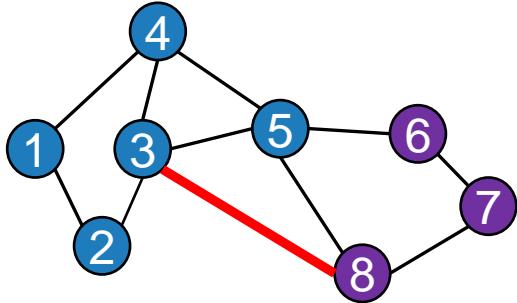
# Adversarial Attacks on GNNs

- Adversarial example on GNNs aims to change the prediction of GNNs by modifying the edges and features.

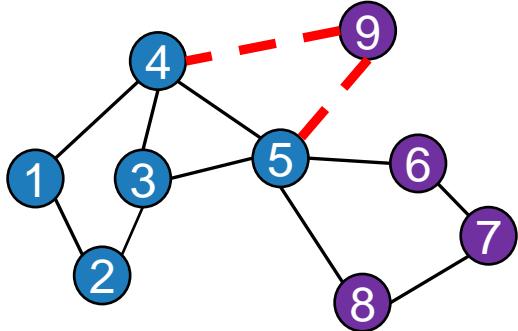


# Types of Perturbations

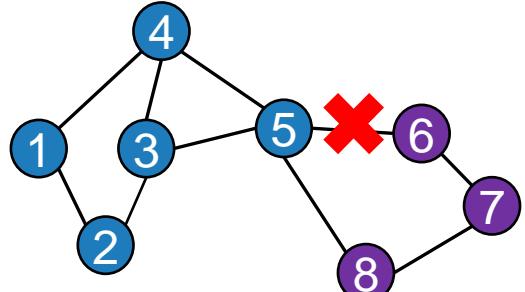
**Edge Insertion**



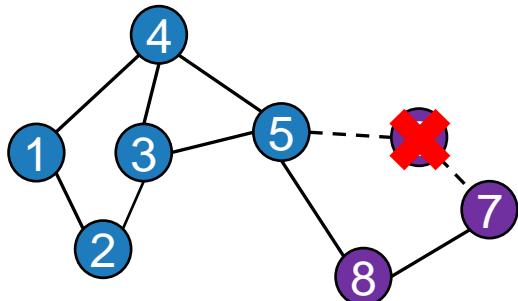
**Node Injection**



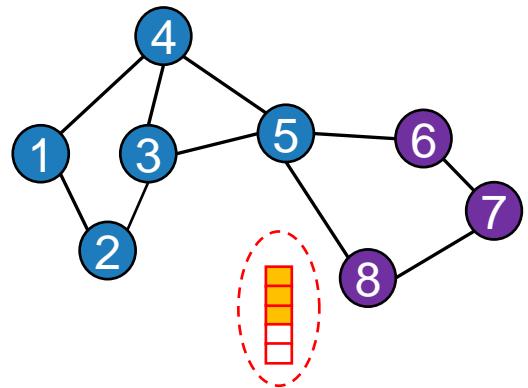
**Edge Deletion**



**Node Deletion**

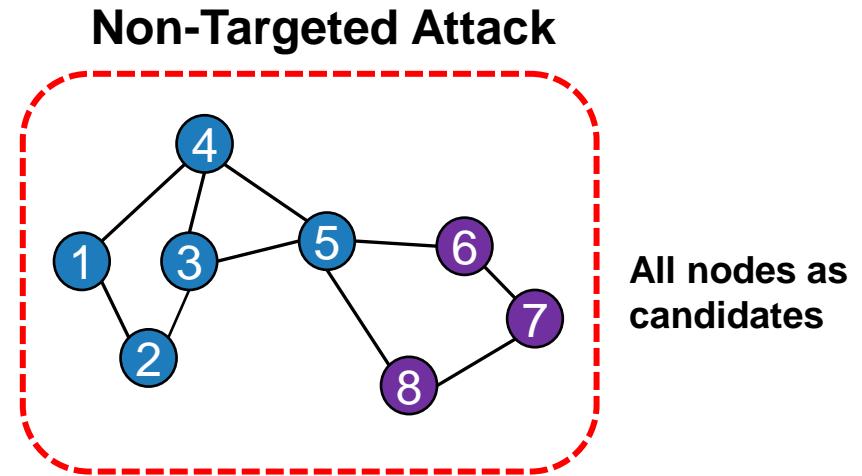
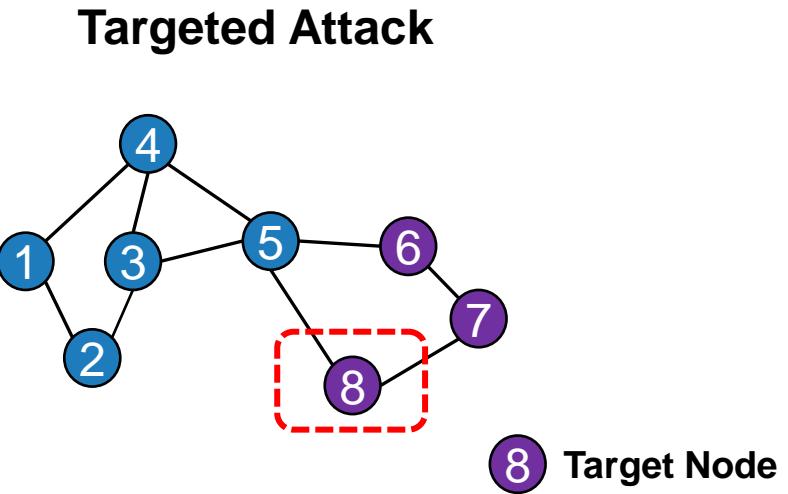


**Feature Modification**

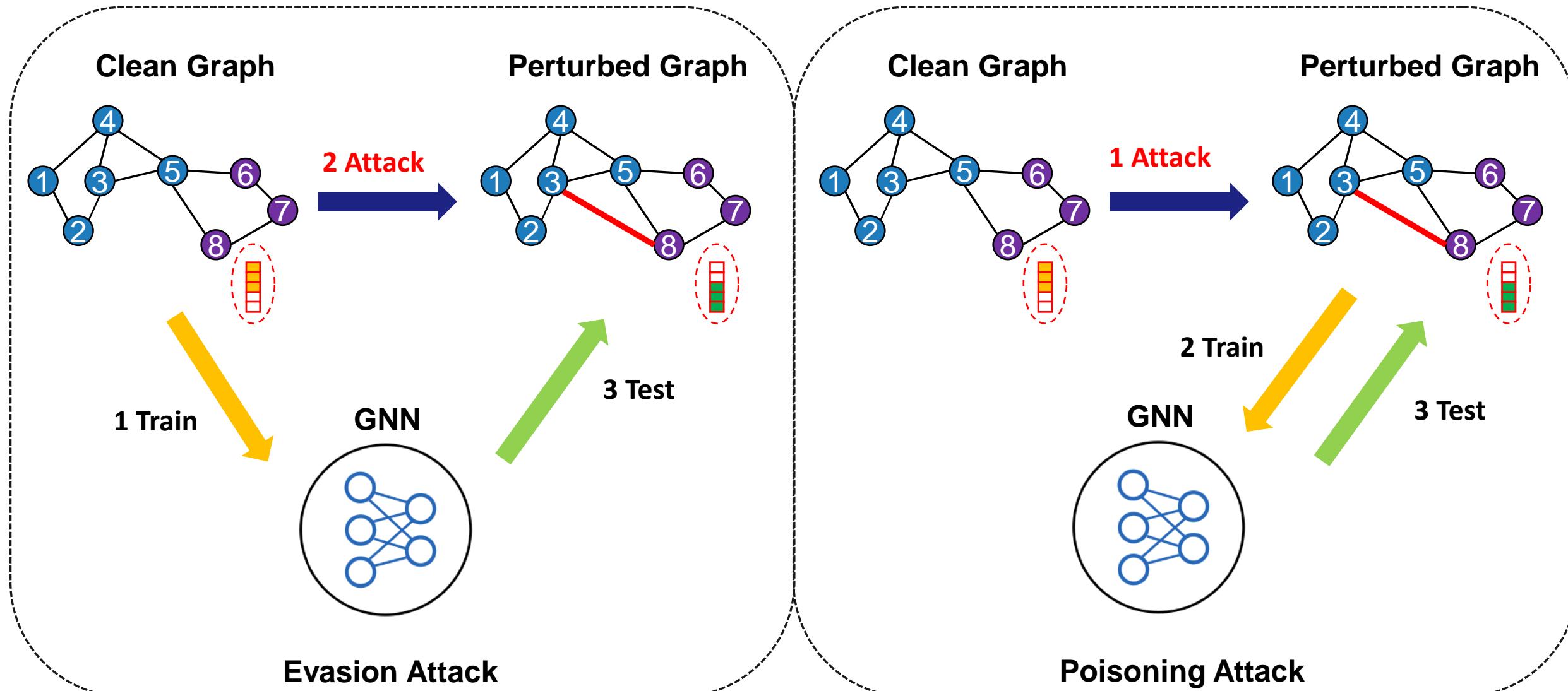




# Types of Attacking Type



# Evasion & Poisoning Attack



# Attacker's Knowledge

- Attacker's knowledge means how much information an attacker knows about the model that he aims to attack.

Settings	Model Parameters	Model Predictions	Labels	Training Input
<b>White-Box Attack (WBA)</b>	✓	✓	✓	✓
<b>Gray-Box Attack (GBA)</b>		✓	✓	✓
<b>Restrict Black-box Attack (RBA)</b>				✓

# Attack Methods

Attack Methods	Node Injection	Edge Insertion/ Deletion	Feature Modification	Targeted	Non- Targeted	Evasion	Poisoning	Black-box	RL based
Nettack (KDD 18)		✓	✓	✓		✓	✓		
RL-S2V (ICML 18)		✓		✓		✓		✓	✓
$\mathcal{A}_{class}$ (ICML 19)		✓		✓	✓		✓	✓	
Metattack (ICLR 19)		✓			✓		✓		
GF-Attack (AAAI 20)		✓		✓		✓		✓	
CD-Attack (WWW 20)		✓		✓			✓	✓	
NIPA (WWW 20)	✓				✓		✓		✓

# Attack Methods

Attack Methods	Node Injection	Edge Insertion/ Deletion	Feature Modification	Targeted	Non- Targeted	Evasion	Poisoning	Black-box	RL based
Nettack (KDD 18)		✓	✓	✓		✓	✓		
RL-S2V (ICML 18)		✓		✓		✓		✓	✓
$\mathcal{A}_{class}$ (ICML 19)		✓		✓	✓		✓	✓	
Metattack (ICLR 19)		✓			✓		✓		
GF-Attack (AAAI 20)		✓		✓		✓		✓	
CD-Attack (WWW 20)		✓		✓			✓	✓	
NIPA (WWW 20)	✓				✓		✓		✓

# Nettack (KDD 18)

- General form of graph adversarial attack as a bi-level optimization problem (poisoning setting):

$$\arg \max_{\hat{A}, \hat{X}} \mathcal{L}_{atk}(f_{\theta^*}(A', X')) = \sum_{u \in V_t} \ell(f_{\theta^*}(A', X')_u, c_{old,u})$$

Prediction of attacked model

Prediction of clean model

s.t.  $\theta^* = \arg \min_{\theta} \mathcal{L}_{train}(f_{\theta}(A', X'))$ ,  $|A' - A| + |X' - X| \leq \Delta$

- Types of attack:**

- Targeted attack, Poison/Evasion Attack, White-Box Attack
- Structure and feature modifications

- Core idea:** Establishing a linear **surrogate model**:

$$f_{S,\theta}(A, X) = \text{softmax}(\hat{A} \text{ReLU}(\cancel{X\Theta^{(1)}})\Theta^{(2)}) = \text{softmax}(\hat{A}^2 X \Theta)$$

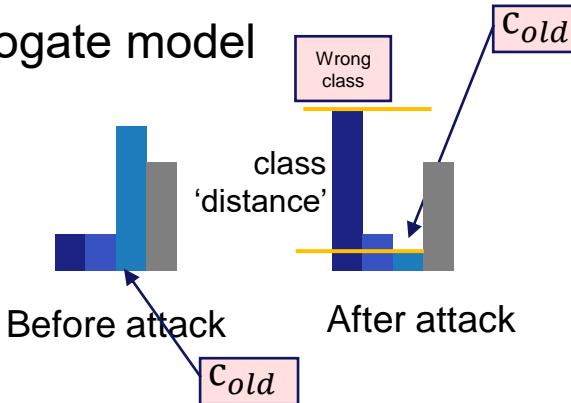
$A$ : adjacent matrix  
 $X$ : node features matrix  
 $A'$ : modified structure  
 $X'$ : modified feature  
 $V_t$ : set of target nodes  
 $c_{old,u}$ : the predicted class label of the clean model.  
 $\Delta$ : perturbation budget

$\Theta^{(1)}\Theta^{(2)}$ , trained on the clean data.

# Nettack (KDD 18)

**Optimization approach:** one perturbation that maximizes the class ‘distance’ of the surrogate model  $f_{S,\theta}(A, X)$  before and after attack.

$$\ell(f_{S,\theta^*}(A, X), c_{old,u}) = \max_{c \neq c_{old}} Z(u)_c - Z(u)_{c_{old}} \quad Z = f_{S,\theta^*}(A, X)$$

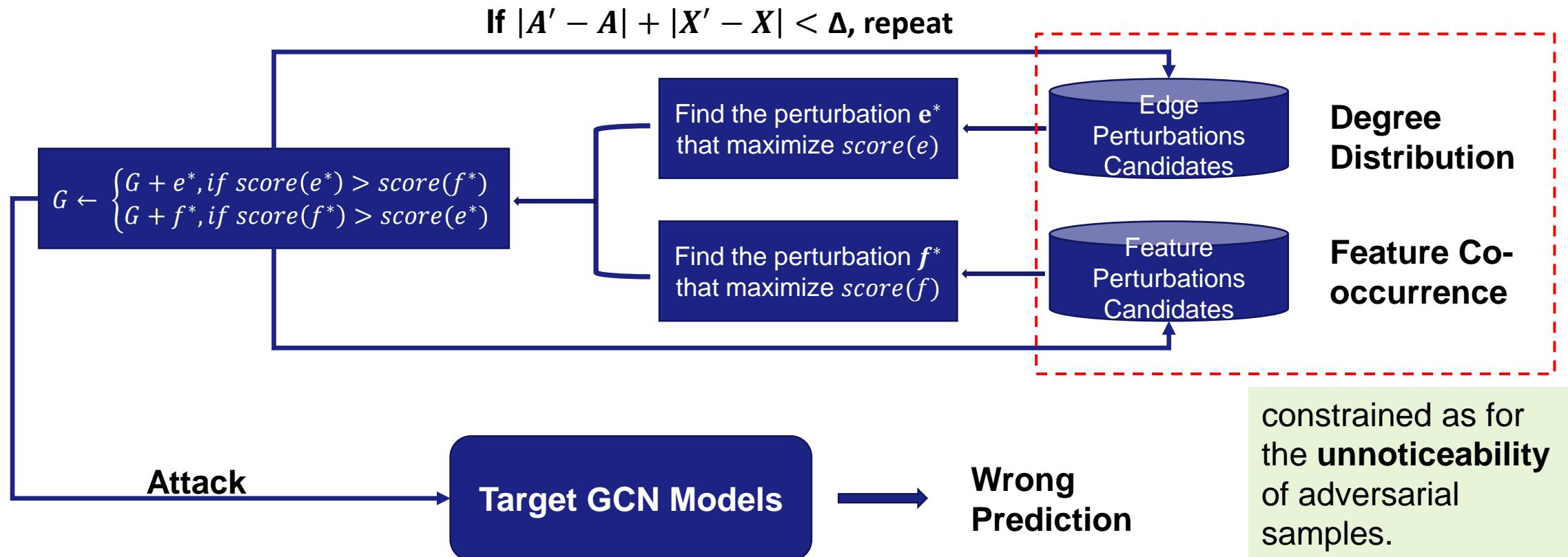


Attack on edge:  $score(e) = \ell(f_{S,\theta^*}(A', X), c_{old,u})$ ,  $A' := A \pm e$

Attack on feature:  $score(f) = \ell(f_{S,\theta^*}(A, X'), c_{old,u})$ ,  $X' := X \pm f$

# Nettack (KDD 18)

**Optimization approach:** one perturbation that maximizes the class ‘distance’ of the surrogate model  $f_{S,\theta}(A, X)$  before and after attack.



## Overview

# GF-Attack (AAAI 20)

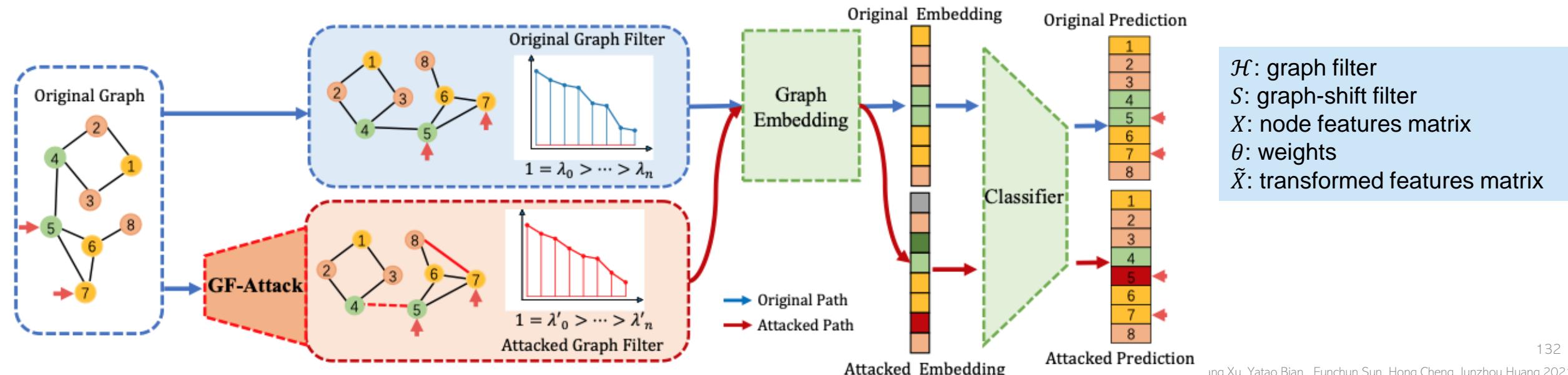
- How to perform the attack under RBA setting?
- View of graph signal processing → attacking on the graph filter(**GF-Attack**).

Graph filtering:  $\tilde{X} = \mathcal{H}(X)$

Feature transformation:  $X' = \sigma(\tilde{X}\Theta)$

- **Types of attack:**

- Target attack, Restricted Black-Box Attack
- Structure modifications



# GF-Attack (AAAI 20)

- **Core idea:** attacking the graph filter  $\mathcal{H}_{\text{model}} = h(S)$  for different graph embedding models.
- Measuring the quality of output embedding  $Z$  as the  $T$ -rank approximation problem:

$$\mathcal{L}(A', Z) = \|h(S')X - h(S')_T X\|_F^2,$$

- By solving this problem, this is equivalent to optimize:

$$\begin{aligned} \arg \max_{A'} \sum_{i=T+1}^n |\lambda'_i|^2 \cdot \sum_{i=T+1}^n \|\mathbf{u}_i^T X\|_2^2, \\ \text{s.t. } \|A' - A\| = 2\beta. \end{aligned}$$

$\mathcal{H}$ : graph filter  
 $S$ : graph-shift filter  
 $h(\cdot)$ : a polynomial function.  
 $X$ : node features matrix  
 $\theta$ : weights  
 $\tilde{X}$ : transformed features matrix  
 $\lambda_i$  and  $\mathbf{u}_i$  are an eigen-pair of graph filter  $\mathcal{H}$ .

Eigen-pair of graph filter  $h(S)$

If we know the graph filter of target model, we can perform the adversarial attack without any additional information(model parameters, predictions, labels)!

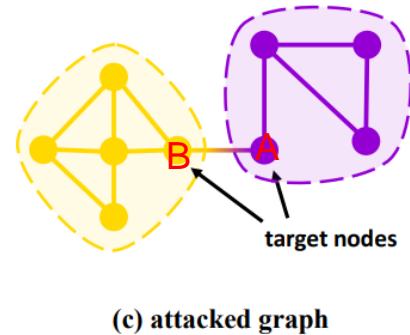
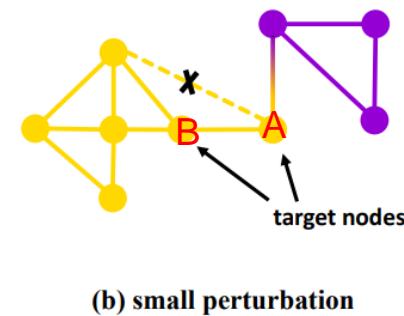
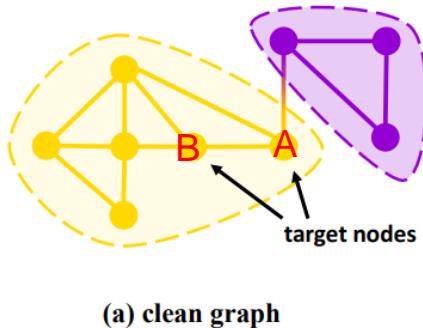
# GF-Attack (AAAI 20)

Graph Embedding Models	Graph-shift filter $S$	Polynomial Function $h(x)$
GCN	$L^{sym} - I_n$	$h(x) = x$
SGC	$L^{sym} - I_n$	$h(x) = x$
ChebyNet	$L^{sym} - I_n$	$h(x) = \sum_{k=0}^K T_k(x)$
LINE	$I_n - L^{rw}$	$h(x) = x$
DeepWalk	$I_n - L^{rw}$	$h(x) = \sum_{k=0}^K x^k$

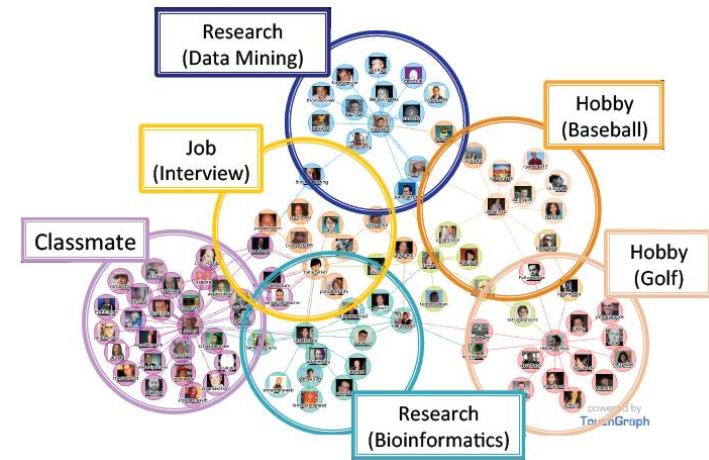
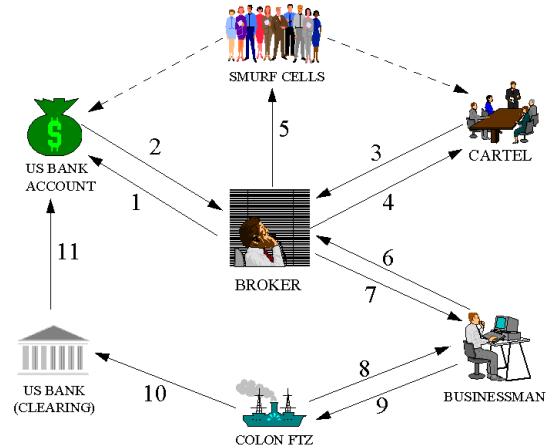
$S$ : graph-shift filter;  
 $h(x)$ : polynomial function used for constructing graph filter  $\mathcal{H}$ .

# CD-Attack (WWW 20)

- Go beyond node classification → hiding individuals from **community detection**.
- **Broad applications:**
  - Money laundering in transaction networks.
  - Personal privacy protection in social networks.
- **Types of attack:**
  - Target attack, Restricted Black-box Attack
  - Structure modifications

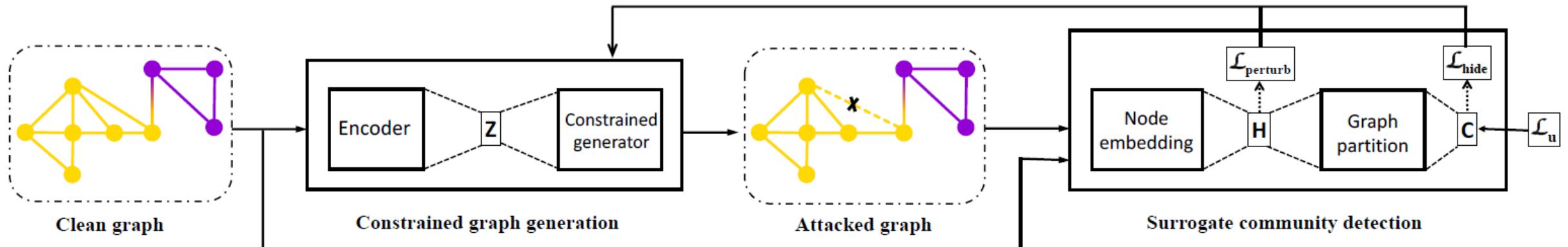


**Target:** hiding individuals from community detection by small modification on the graph structure.



# CD-Attack (WWW 20)

- **Challenges:**
  - Huge space to locate the perturbation positions
  - No knowledge about the targeted community detection methods (RBA)
- **Core idea:** Modeling the interaction between community detection and perturbed graph:
  - **Constrained graph generation:** VGAE with efficient sampling.
  - **Surrogate community detection:** GNN (efficiency) and normalized cut (generalization)
  - **Interaction between two modules:** Actor-Critic

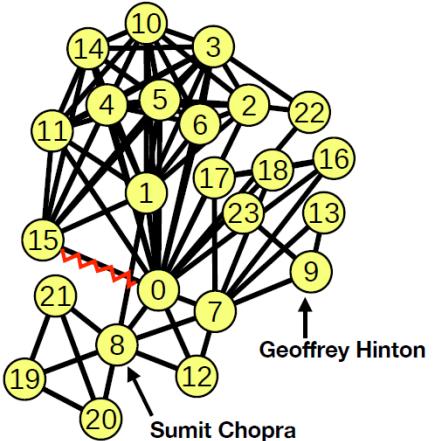


# Experimental Results.

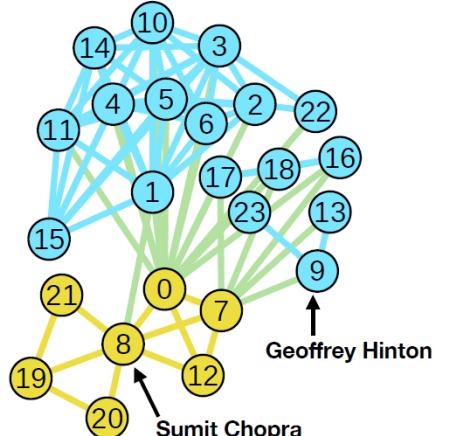
Data sets	DBLP-medium		Finance-medium		DBLP-large		Finance-large	
	M1	M2	M1	M2	M1	M2	M1	M2
DICE	6.35%	41.44%	3.25%	53.55%	6.91%	50.77%	3.66%	72.11%
MBA	7.33%	44.88%	2.80.%	52.00%	7.04%	53.67%	3.08%	72.11%
RTA	6.55%	42.22%	2.53%	42.66%	6.25%	44.66%	1.74%	55.11%
CD-ATTACK	13.72%	52.00%	3.32%	63.00%	8.11%	57.62%	4.14%	87.33%

The hiding quality  
(Higher the better)

Case Study: Attack node2vec (Grover and Leskovec, 2016) + K-means on DBLP.



(a) clean graph



(b) modified graph

- 0: Yann LeCun
- 1: Raia T. Hadsell
- 2: Koray Kavukcuoglu
- 3: Pierre Sermanet
- 4: Ayse Erkan
- 5: Jan Ben
- 6: Urs Muller
- 7: M. Ranzato
- 8: Sumit Chopra
- 9: Geoffrey Hinton
- 10: Marco Scoffier
- 11: Matt Grimes
- 12: Fu Jie Huang
- 13: Volodymyr Mnih
- 14: Urs Miller
- 15: Chris Crudele
- 16: J. M. House
- 17: R. C. Flagan
- 18: P. E. Taylor
- 19: Andrew Caplin
- 20: John Leahy
- 21: Trivikraman Thampy
- 22: Soumith Chintala
- 23: Andrew Denis Brown

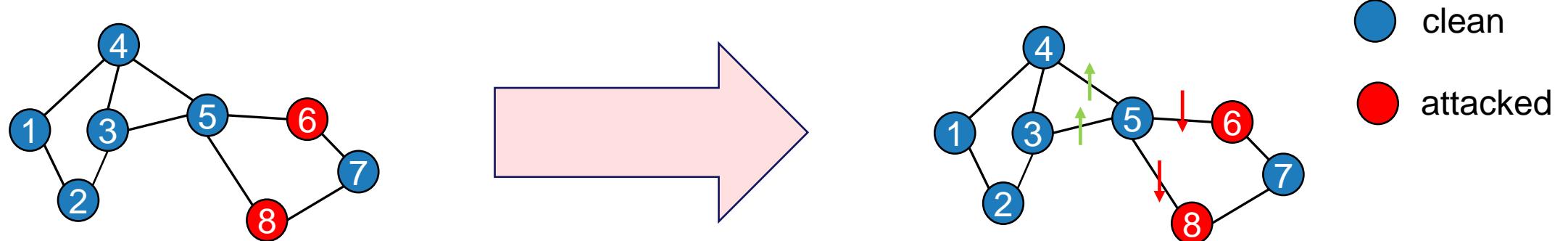
# Defense against Adversarial Attacks

Defense strategies on graph data can be categorized as:

- **Attention mechanism**
- **Graph purification (preprocessing)**

# Attention mechanism (RobustGCN)

- Using the attention mechanism to reduce the influence of adversarial attack.



- RobustGCN** assumes that the adversarial sample nodes may have **high prediction uncertainty**.
- Core idea:** Giving lower attention scores to adversarial edges to reduce their impact.

# Attention mechanism (RobustGCN)

- Prediction uncertainty is captured by **Gaussian-based Graph Convolution Layer** (GGCL).

- Approach:**

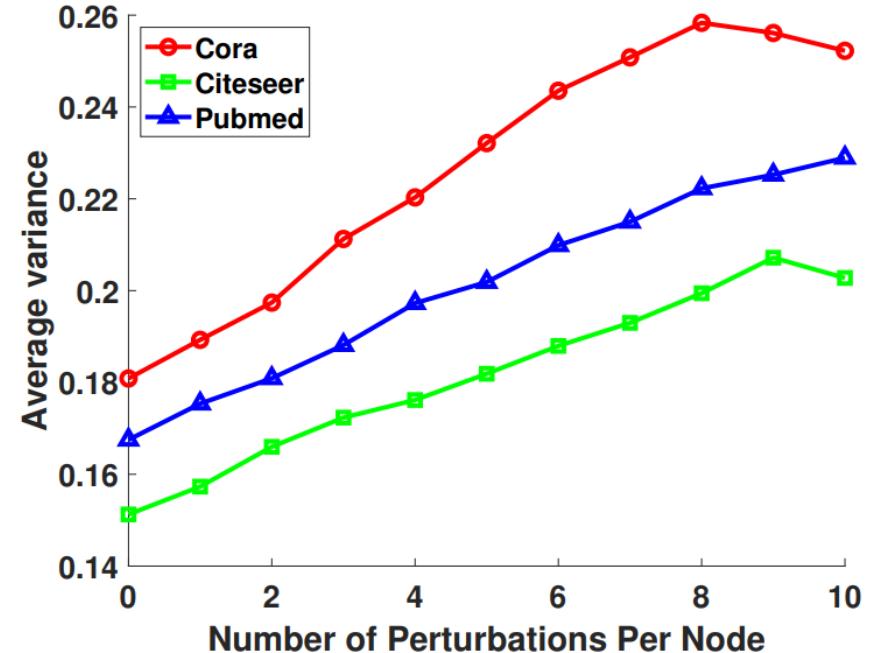
- Using Gaussian distribution to capture the uncertainty in GGCL:

$$\mathbf{h}_{N(i)}^{(l)} = \sum_{j \in N(i)} \frac{1}{\sqrt{\tilde{\mathbf{D}}_{ii} \tilde{\mathbf{D}}_{jj}}} \mathbf{h}_j^{(l)} \sim N\left(\sum_{j \in N(i)} \frac{1}{\sqrt{\tilde{\mathbf{D}}_{ii} \tilde{\mathbf{D}}_{jj}}} \boldsymbol{\mu}_j^{(l)}, \text{diag}\left(\sum_{j \in N(i)} \frac{1}{\tilde{\mathbf{D}}_{ii} \tilde{\mathbf{D}}_{jj}} \sigma_j^{(l)}\right)\right)$$

- Applying the attention mechanism to assign the nodes with **high variance** with lower weights when aggregation:

$$\mathbf{h}_{N(i)}^{(l)} = \sum_{j \in N(i)} \frac{1}{\sqrt{\tilde{\mathbf{D}}_{ii} \tilde{\mathbf{D}}_{jj}}} (\mathbf{h}_j^{(l)} \odot \alpha_j^{(l)})$$

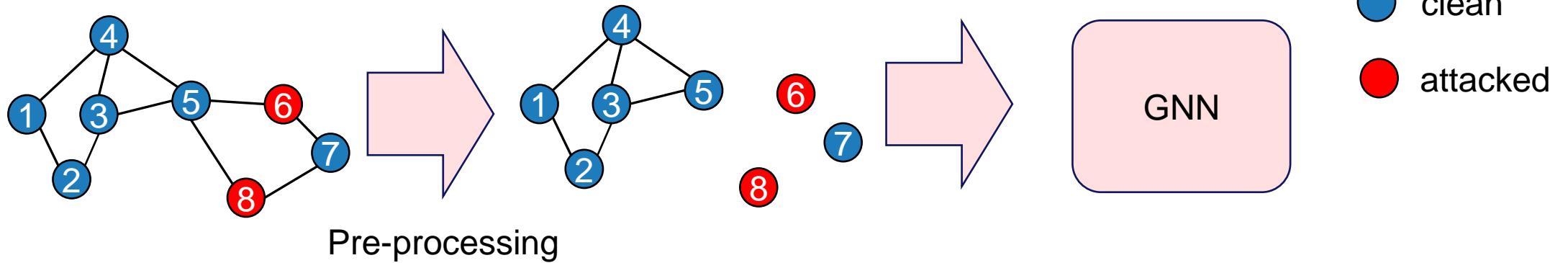
$$\alpha_j^{(l)} = \exp(-\gamma \sigma_j^{(l)})$$



high variance → more perturbation

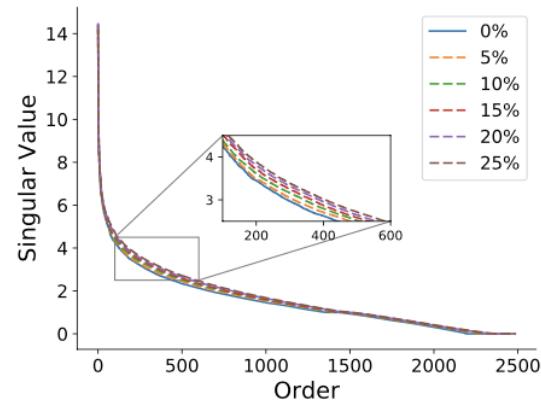
# Graph purification (Pro-GNN)

- The purification method firstly purifies the disturbed graph structure data, and then trains the GNN model on the purified graph.

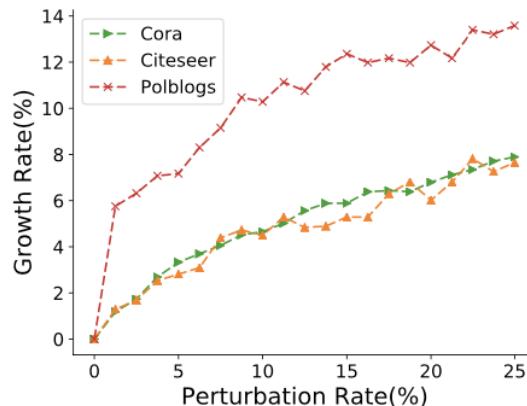


# Graph purification (Pro-GNN)

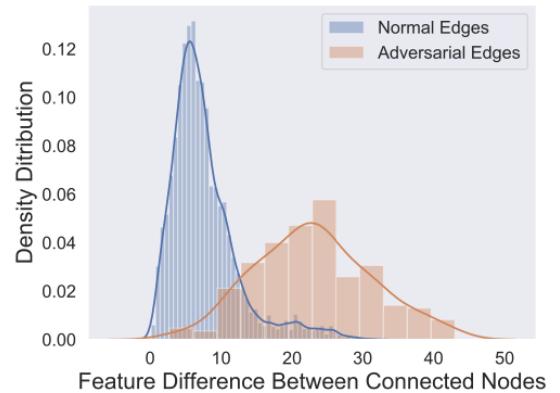
- **Pro-GNN:** explores the intrinsic properties that shared by real-world graph and uses these properties to purify poisoned graphs as well as learn a robust GNN.



Larger singular values



Larger rank



Non-smooth feature

# Graph purification (Pro-GNN)

- Approach:

- Add loss to ensure the graph is **low-rank** and **sparse**:

$$\arg \min_{S \in S} \mathcal{L}_0 = \|A - S\|_F^2 + \alpha \|S\|_1 + \beta \|S\|_*, \text{ s.t., } S = S^\top,$$

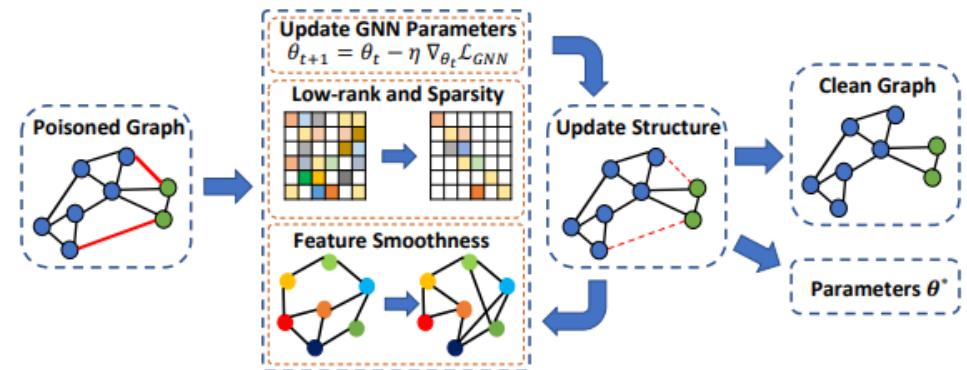
- Add loss to penalize the **rapid changes in features** between adjacent nodes:

$$\arg \min_{S \in S} \mathcal{L} = \mathcal{L}_0 + \lambda \cdot \mathcal{L}_s = \mathcal{L}_0 + \lambda \operatorname{tr}(X^T \hat{L} X), \text{ s.t., } S = S^\top,$$

- Jointly learning the desired properties of graphs and the GNN model:

$$\begin{aligned} \arg \min_{S \in S, \theta} \mathcal{L} &= \mathcal{L}_0 + \lambda \mathcal{L}_s + \gamma \mathcal{L}_{GNN} \\ &= \|A - S\|_F^2 + \alpha \|S\|_1 + \beta \|S\|_* + \gamma \mathcal{L}_{GNN}(\theta, S, X, Y_L) + \lambda \operatorname{tr}(X^T \hat{L} X) \\ \text{s.t. } & S = S^\top, \end{aligned}$$

$S$ : the clean adjacent matrix we would like to learn



Effective but time-consuming.



# Defense Methods

Defense Methods	Attention Mechanism	Preprocessing
GNN-Jaccard (IJCAI 19)		✓
<b>RGCN</b> (KDD 19)	✓	
<b>Pro-GNN</b> (KDD 20)		✓
GNN-SVD (WSDM 20)		✓
<b>GNNGUARD</b> (NeurIPS 20)	✓	
VPN (AAAI 21)	✓	



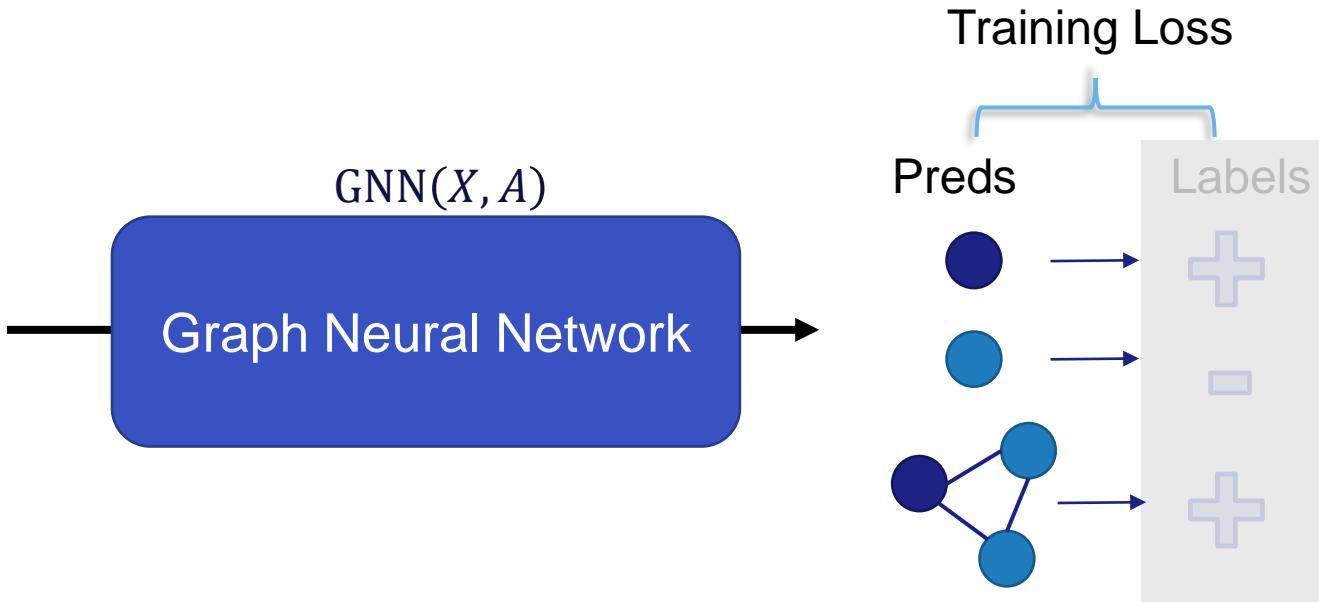
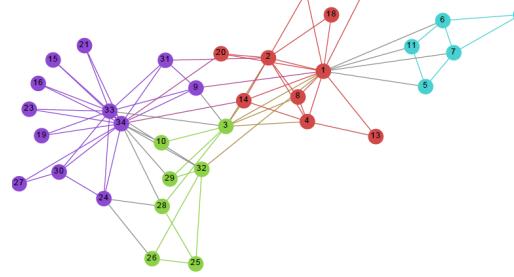
Tencent  
AI Lab



# **Self/Un-Supervised Learning of GNNs**



# What we discussed before are supervised



- **Labels are scarce**, e.g. molecular property
- **Training/Testing tasks are Non I.I.D.**

# Existing Self-Supervised GNNs

	Node- Classification	Link/Metapath Prediction	Graph- Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S <sup>2</sup> GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5] , GROVER [7] GCC [6]	
Information- based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]



- [1] Kipf & Welling 2016; [2] Durán & Niepert 2017; [3] Hamilton et al. 2017;
- [4] Liu et al. 2019; [5] Hu et al. 2020; [6] Qiu et al. 2020; [7] Rong et al. 2020;
- [8] Veličković et al. 2019; [9] Peng et al. 2020; [10] Sun et al. 2020
- [11] Peng, Zhen, et al. 2020 [12] Hwang, Dasol, et al. 2020
- [13] Pan, Shirui, et al. 2018 [14] Hasanzadeh, Arman, et al. 2019



“In self-supervised learning, the system learns to predict part of its input from other parts of its input.” ---- by Yann Lecun

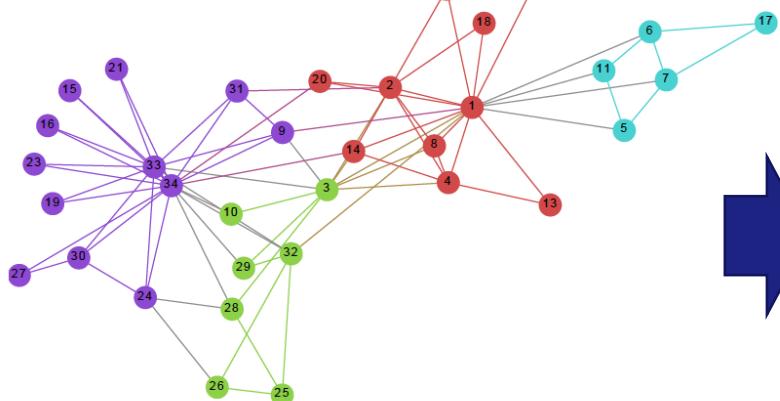
**Graphs are highly structured!**





# Node Classification

- Two typical ways to formulate training loss

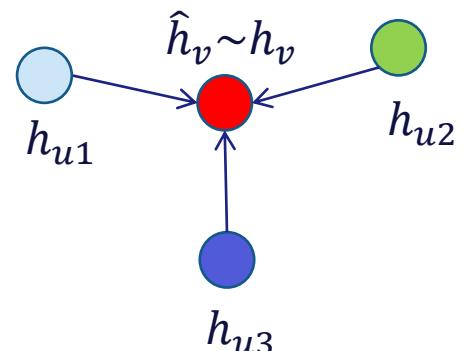
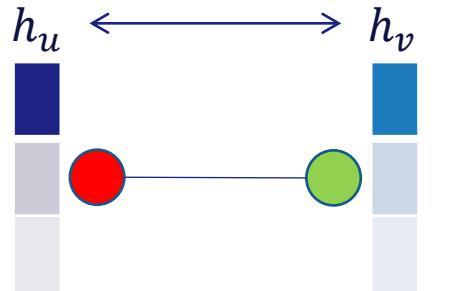


I. Enforcing  
Adjacent Similarity

GraphSAGE

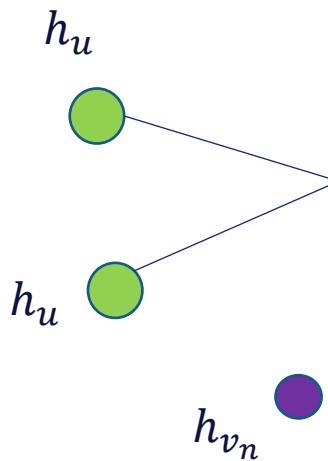
II. Reconstruction  
from Neighbors

EP-B



# I Enforcing Adjacent Similarity

- GraphSAGE (Hamilton et al. 2017)



Enforcing nearby nodes to have similar representations, while enforcing disparate nodes to be distinct:

$$\min -E_{u \sim N(v)} \log(\sigma(h_u^T h_v)) - \lambda E_{v_n \sim P_n(v)} [\log(\sigma(-h_{v_n}^T h_v))]$$

Positive Samples

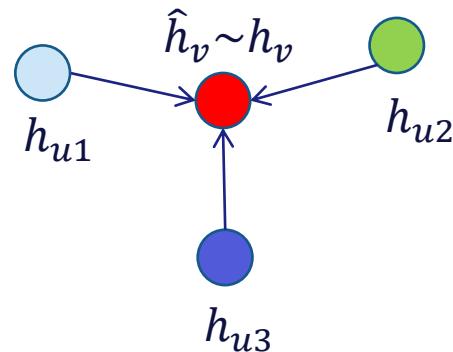
Negative Samples

$h_v$ : representation of **target** node;  
 $h_u$ : representation of **neighbor/positive** node;  
 $h_{v_n}$ : representation of **disparate/negative** node;  
 $P_n(v)$ : negative sampling.

## II Reconstruction from neighbors

- EP-B (Durán & Niepert, 2017)

The objective is to minimize the reconstruction error (regulated by the error to other nodes):



$$\min \sum_{u \in V \setminus \{v\}} [\gamma + d(\tilde{\mathbf{h}}_v, h_v) - d(\tilde{\mathbf{h}}_v, h_u)]_+$$

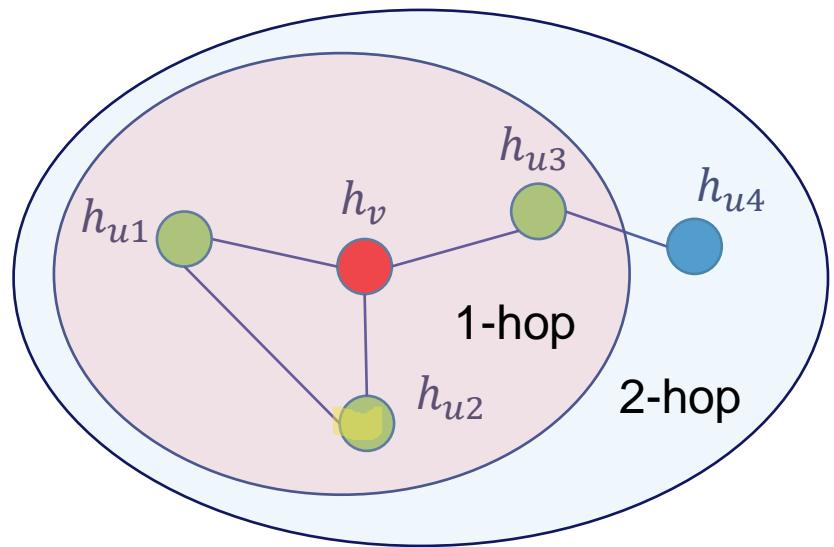
Positive Samples      Negative Samples

Hinge loss

$h_v$ : representation of **target** node;  
 $h_u$ : representation of **nodes except  $v$** ;  
 $\tilde{\mathbf{h}}_v$ :  $\text{AGG}(h_l | l \in N(v))$  is the reconstruction from neighbors;  
 $\gamma$ : the bias

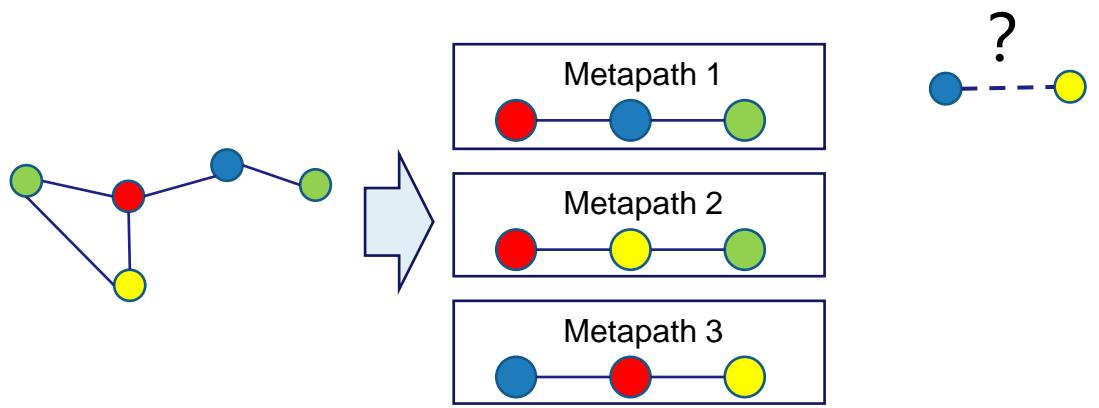
# Link/Metapath Prediction

- S<sup>2</sup>GRL(Peng, Zhen, et al. 2020 )



Predict hop counts (K-hop connectivity)

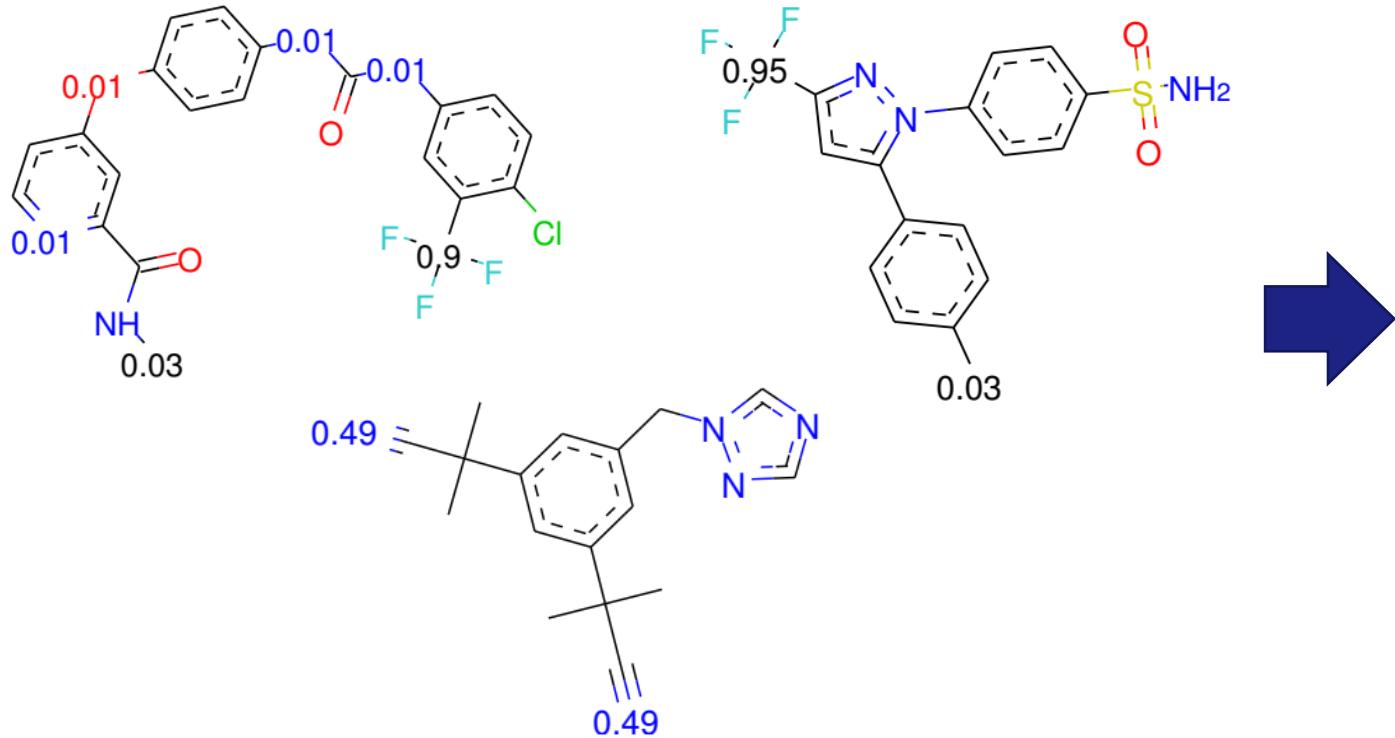
- SELAR(Hwang, Dasol, et al. 2020 )



Predict the type of meta path between two nodes



# How about graph classification/regression?

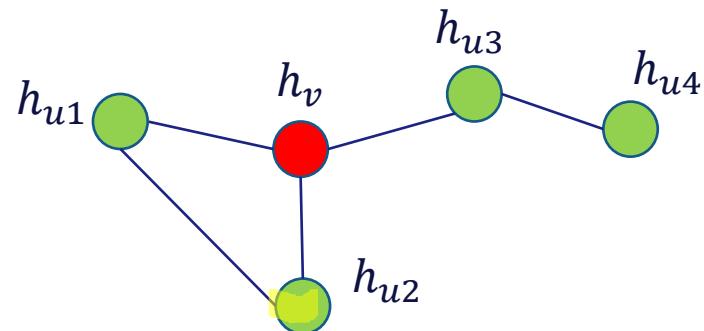


Toxicity?  
Solubility?  
...

# N-Gram Graph

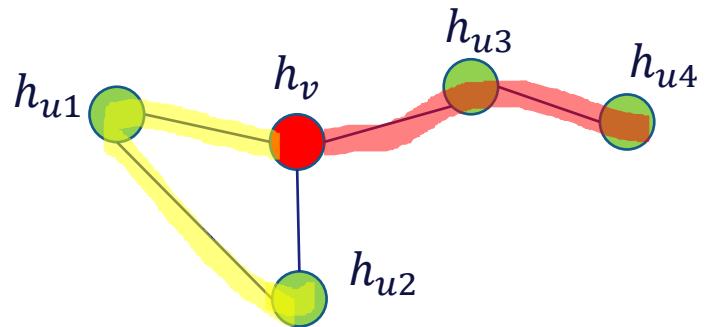
- (Liu et al. 2019)

Stage I: Node Representation



First learn node representations by CBoW-like pipeline

Stage II: Graph Representation



For all n-gram paths:

$$f_p = \prod_{i \in p} h_i ;$$

$$f_{(n)} = \sum_{p \in \text{n-gram}} f_p$$

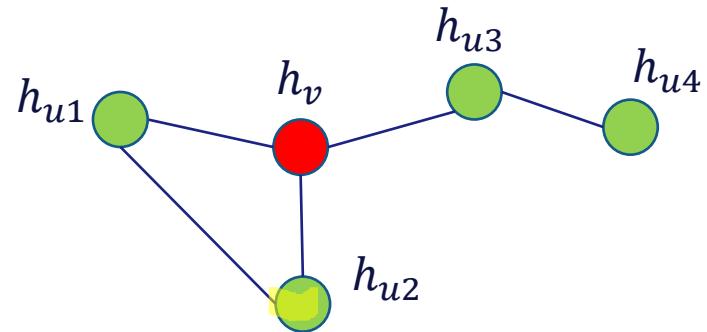
Graph Representation:  $F = [f_{(1)}, \dots, f_{(T)}]$

Equivalent to a GNN that needs no training

# PreGNN: Node- and Graph-level Pretraining

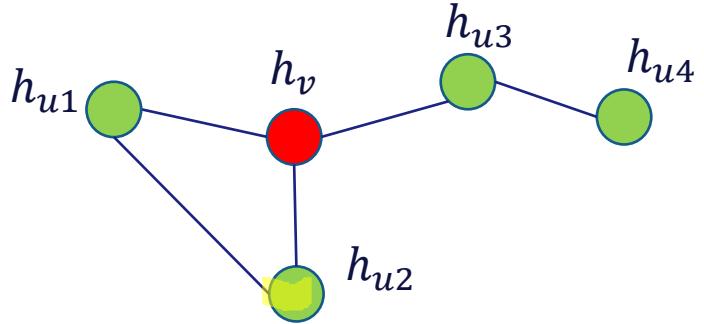
- (Hu et al. 2020)

Stage I: Node Representation



First learn node representations by **Context Prediction or Attribute Masking**

Stage II: Graph Representation



Then perform graph-level multi-task ***Supervised Training***

$$h_G = \text{Readout}(h_v | v \in G)$$

$$\min \text{ CrossEntropy}(h_G, y_G)$$

Both node- and graph- level training are crucial!



# PreGNN

- (Hu et al. 2020)

## Stage I: Node Representation

Enforcing node representation to be similar to its contextual structures:

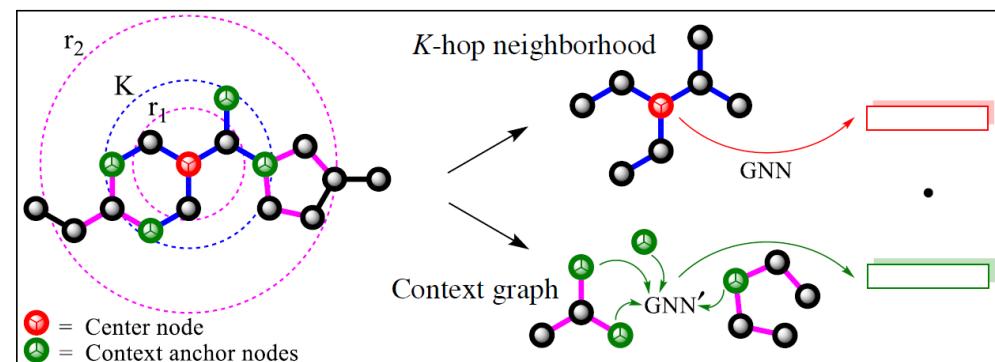
$$\min -\log \left( \sigma((h_v^K)^T C_v^G) \right) - I(v \neq v') \log(1 - \sigma((h_v^K)^T C_{v'}^{G'}))$$

Positive Samples

Negative Samples

Degenerates to EP-B,  
if  $r_1=0, r_2=K=1$

### Context Prediction Attribute Masking



$h_v^K$ : K-hop information

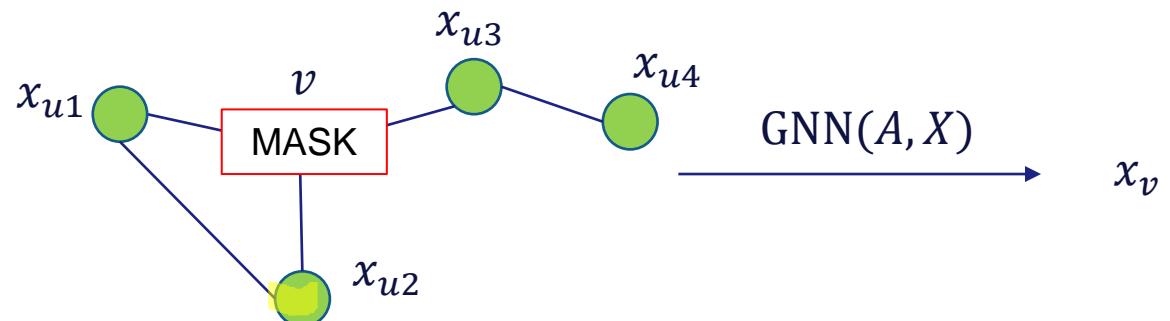
$C_v^G$ : Structures between  
 $r_1$  and  $r_2$ -hop

# PreGNN

- (Hu et al. 2020)

Stage I: Node Representation

Mask random node/edge attribute and predict it, just like Bert:

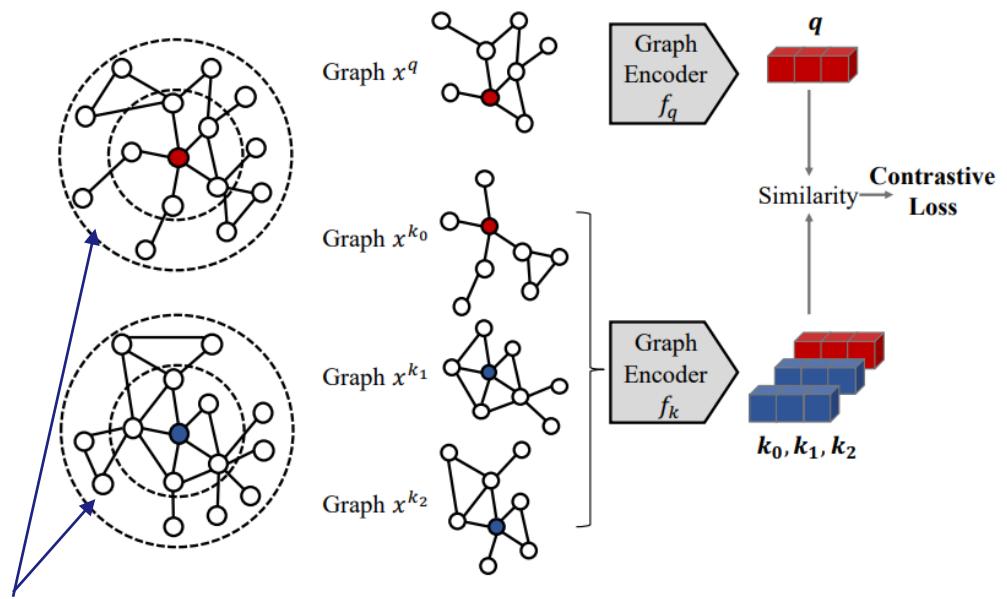


**Context Prediction  
Attribute Masking**

# GCC: Contrastive learning

- (Qiu et al. 2020)

Both N-Gram Graph and PreGNN do **not** perform **graph-level unsupervised** training:



r-ego network

$$\min - \log \frac{\exp(q^T k_+ / \tau)}{\sum_{i=0}^K \exp(q^T k_i / \tau)}$$

InfoNCE

$q$ : representation of different graphs;  
 $k$ : key of different graphs;  
 $k_+$ : positive key generated by random graph perturbation  
 $\tau$ : temperature

But, GCC only conducts graph-level pre-training, without node-level distinction



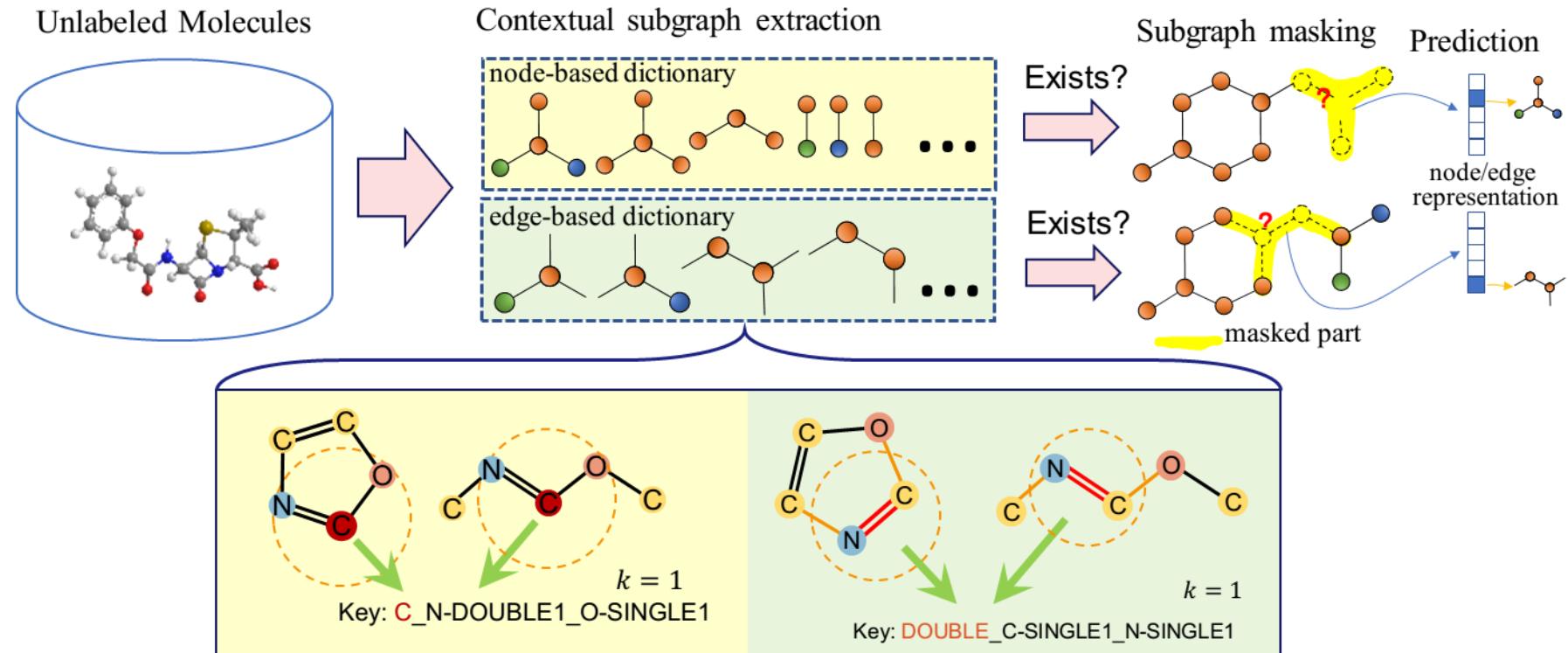
# GROVER (Rong et al. 2020)

Methods	Node-Level Self-Supervised	Graph-Level Self-Supervised
N-Gram Graph	✓	✗
PreGNN	✓	✗
GCC	✗	✓
GROVER	✓	✓

# GROVER

- (Rong et al. 2020)

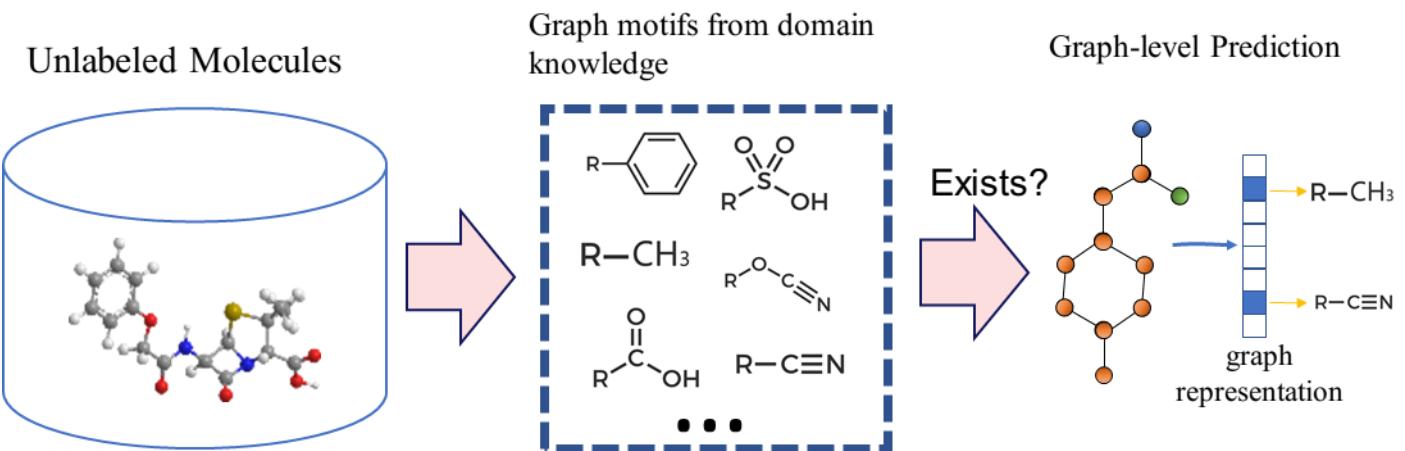
## Stage I: Node/Edge-level pretraining



Predicting node/edge contexts instead of node labels can better capture local structures (multi-label)

- (Rong et al. 2020)

## Stage II: Graph-level pretraining

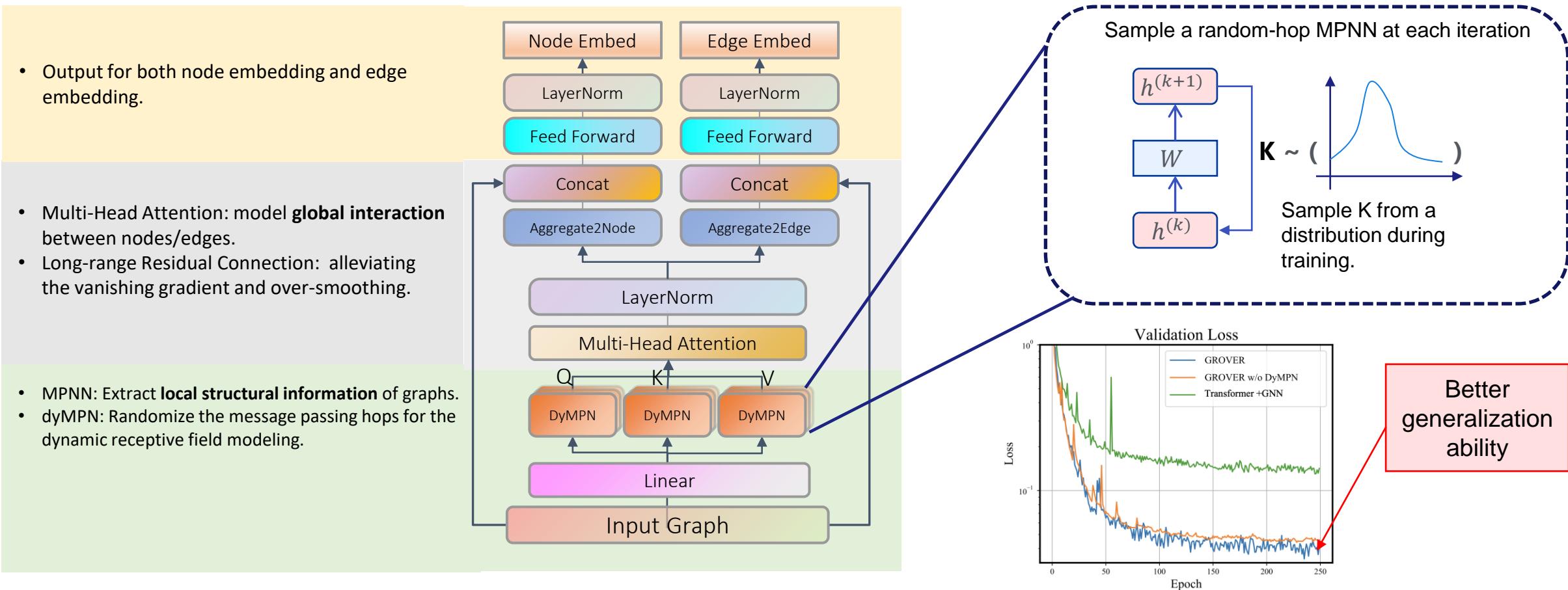


Predicting a graph if contains pre-defined graph motifs.

# GROVER

- One more thing: GTransformer

We build a more expressive and transformer-alike model: GTransformer





We pre-train GROVER with **100 million** parameters on  
**10 million** unlabeled molecules collected from ZINC15  
and ChembI



# GROVER

## Molecular classification

Dataset # Molecules	Classification (Higher is better)					
	BBBP 2039	SIDER 1427	ClinTox 1478	BACE 1513	Tox21 7831	ToxCast 8575
TF_Robust [39]	0.860(0.087)	0.607(0.033)	0.765(0.085)	0.824(0.022)	0.698(0.012)	0.585(0.031)
GraphConv [23]	0.877(0.036)	0.593(0.035)	0.845(0.051)	0.854(0.011)	0.772(0.041)	0.650(0.025)
Weave [22]	0.837(0.065)	0.543(0.034)	0.823(0.023)	0.791(0.008)	0.741(0.044)	0.678(0.024)
SchNet [44]	0.847(0.024)	0.545(0.038)	0.717(0.042)	0.750(0.033)	0.767(0.025)	0.679(0.021)
MPNN [13]	0.913(0.041)	0.595(0.030)	0.879(0.054)	0.815(0.044)	0.808(0.024)	0.691(0.013)
DMPNN [61]	0.919(0.030)	0.632(0.023)	0.897(0.040)	0.852(0.053)	0.826(0.023)	0.718(0.011)
MGCN [29]	0.850(0.064)	0.552(0.018)	0.634(0.042)	0.734(0.030)	0.707(0.016)	0.663(0.009)
AttentiveFP [59]	0.908(0.050)	0.605(0.060)	0.933(0.020)	0.863(0.015)	0.807(0.020)	0.579(0.001)
N-GRAM [28]	0.912(0.013)	0.632(0.005)	0.855(0.037)	0.876(0.035)	0.769(0.027)	- <sup>2</sup>
HU. et.al[18]	0.915(0.040)	0.614(0.006)	0.762(0.058)	0.851(0.027)	0.811(0.015)	0.714(0.019)
GROVER <sub>base</sub>	0.936(0.008)	0.656(0.06)	0.925(0.013)	0.878(0.016)	0.819(0.020)	0.723(0.010)
GROVER <sub>large</sub>	<b>0.940</b> (0.019)	<b>0.658</b> (0.023)	<b>0.944</b> (0.021)	<b>0.894</b> (0.028)	<b>0.831</b> (0.025)	<b>0.737</b> (0.010)

# Existing Self-Supervised GNNs

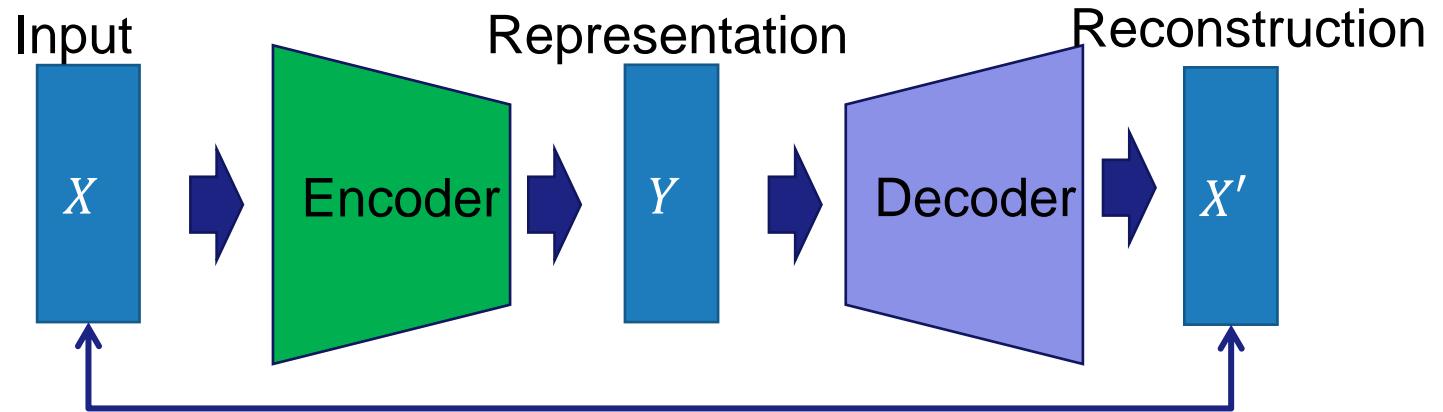
	Node- Classification	Link/Metapath Prediction	Graph- Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S <sup>2</sup> GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5] , GROVER [7] GCC [6]	
Information- based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]



- [1] Kipf & Welling 2016; [2] Durán & Niepert 2017; [3] Hamilton et al. 2017;
- [4] Liu et al. 2019; [5] Hu et al. 2020; [6] Qiu et al. 2020; [7] Rong et al. 2020;
- [8] Veličković et al. 2019; [9] Peng et al. 2020; [10] Sun et al. 2020
- [11] Hwang, Dasol, et al. 2020 [12] Peng, Zhen, et al.
- [13] Pan, Shirui, et al. 2018 [14] Hasanzadeh, Arman, et al. 2019

# What makes a good representation?

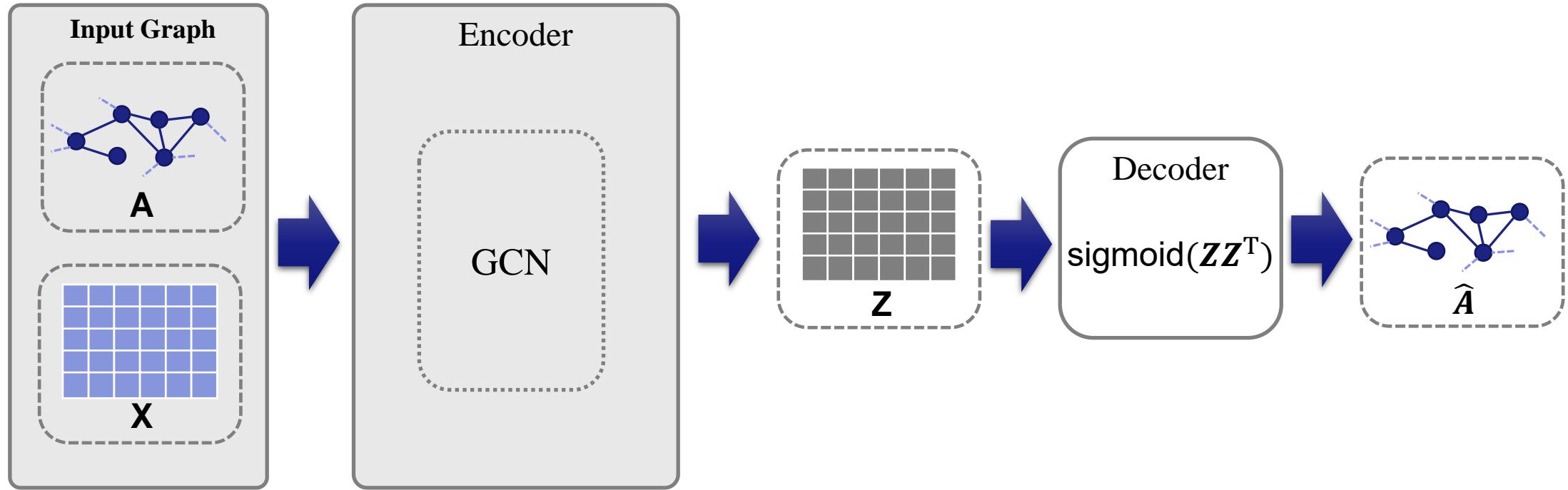
Auto-Encoder (AE)



“One natural criterion that we may expect any good representation to meet, at least to some degree, is to **retain a significant amount of information about the input.**” by Vincent et al. 2010

Hinton & Salakhutdinov 2006; Vincent et al. 2010

# Graph Auto-Encoders (VGAE)



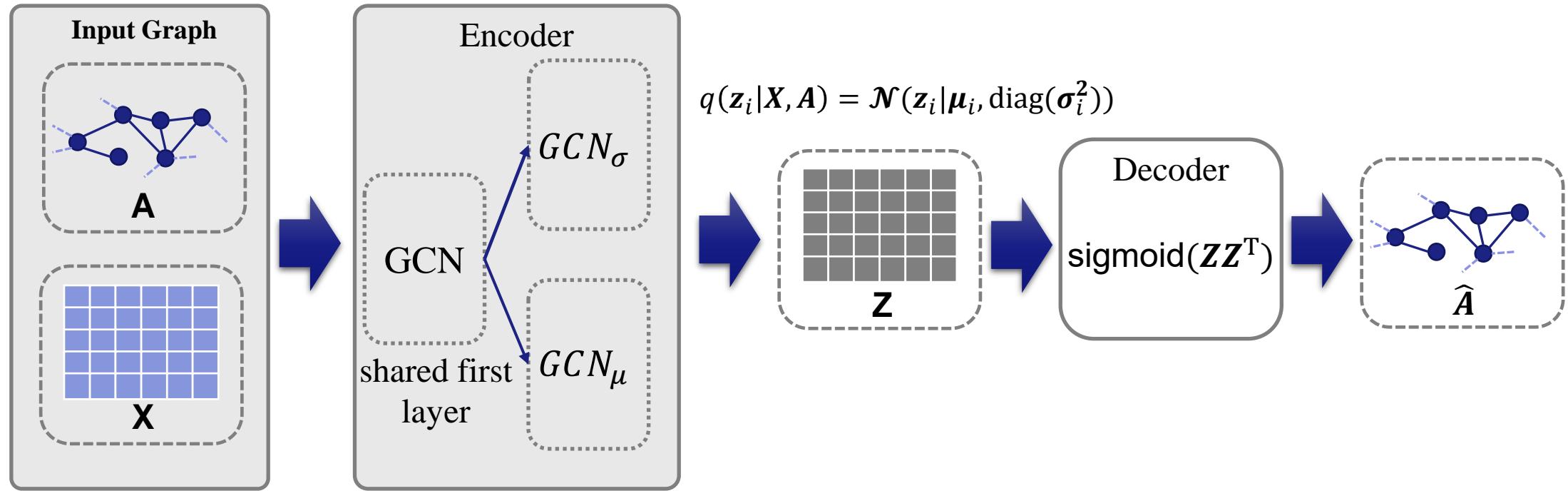
The overall loss:

$$L = \mathbb{E}_{q(\mathbf{Z}|\mathbf{X}, \mathbf{A})} [\log p(\mathbf{A}|\mathbf{Z})]$$

The reconstruction loss.



# Variational Graph Auto-Encoders (VGAЕ)



The overall loss:

$$L = \mathbb{E}_{q(\mathbf{Z}|\mathbf{X}, \mathbf{A})} [\log p(\mathbf{A}|\mathbf{Z})] - \text{KL}[q(\mathbf{Z}|\mathbf{X}, \mathbf{A}) || p(\mathbf{Z})]$$

The reconstruction loss.

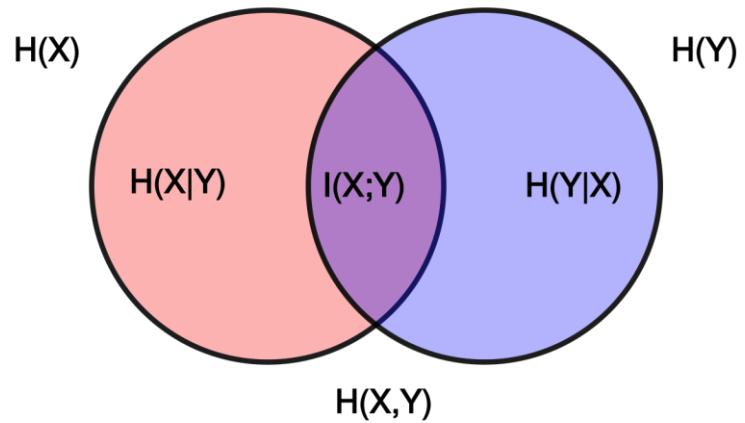
The KL divergence  
between  $q(\cdot)$  and  $p(\cdot)$

# What makes a good representation?

- A more direct way, other than AE?
  - Yes, **Mutual Information (MI)**.

$$\begin{aligned} I(X; Y) &= D_{KL}(p(X)p(Y) \parallel p(X, Y)) \\ &= H(X) - H(X|Y) \end{aligned}$$

Entropy    Conditional Entropy



- $0 \leq I(X; Y) \leq H(X)$  or  $H(Y)$ ;
- $I(X; Y) = 0$  iff  $X$  and  $Y$  are independent random variables;
- $I(X; Y) = H(X) = H(Y)$ , if  $X$  and  $Y$  are determinately related, i.e.  $H(X|Y) = 0$

# AE is a lower bound of MI

(Hjelm et al. 2019)

$$I(X; Y) = H(X) - H(X|Y) \geq H(X) - R(X|Y)$$



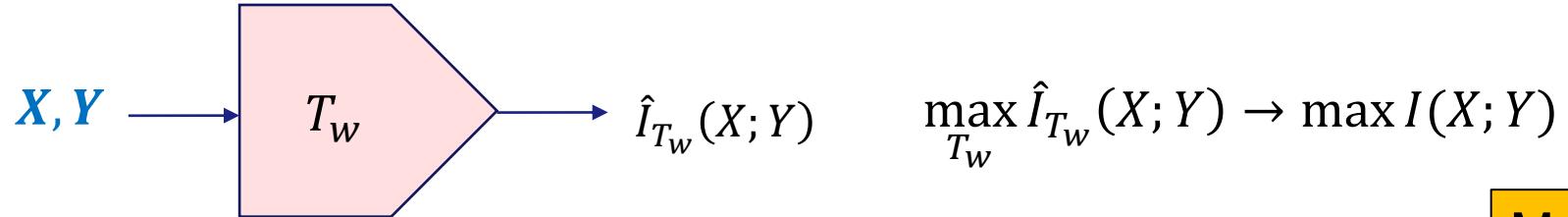
Mutual Information



Reconstruction error

Computing MI is hard and not end-to-end, until recently (CPC, Oord et al., 2018; MINE, Belghazi et al., 2018; Nowozin et al., 2016; Hjelm et al. 2019)

# Estimating/Maximizing MI (Hjelm et al. 2019)



Maximize Lower bound of MI

MINE (Belghazi et al., 2018):

$$I^{\text{MINE}}(X;Y) \triangleq E_{p(X,Y)}[T_w(x,y)] - \log E_{p(X)p(Y)}[\exp(T_w(x,y))]$$

JSD MI estimator (Nowozin et al., 2016):

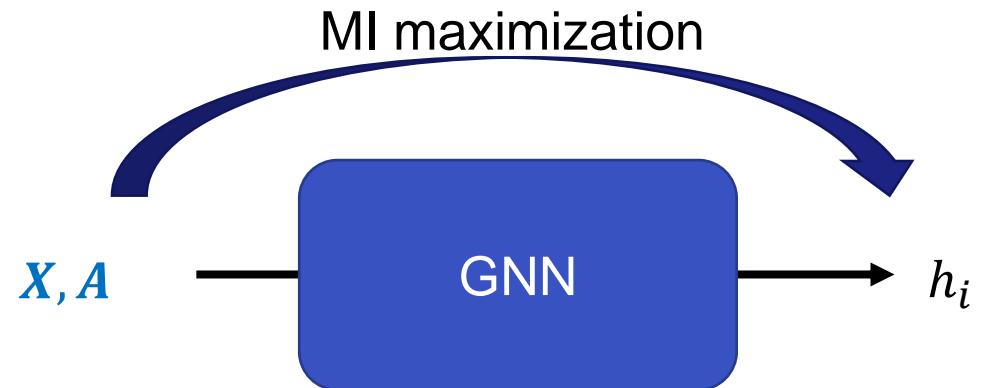
$$I^{\text{JSD}}(X;Y) \triangleq E_{p(X,Y)}[\log \sigma(T_w(x,y))] + E_{p(X)p(Y)}[\log(1 - \sigma(T_w(x,y)))]$$

InfoNCE MI estimator (Oord et al., 2018):

$$I^{\text{NCE}}(X;Y) \triangleq E_{p(X,Y)}[\log \frac{\exp T_w(x,y)}{\sum_{x' \sim p(X)} \exp T_w(x',y)}]$$

# Deep Graph Infomax (DGI)

- (Velickovic et al. 2019)



The JSD MI estimator is applied:

$$\max_{\text{GNN}} I(X, A; h_i) \approx \max \log(D(h_i; \mathbf{X}, \mathbf{A})) + \log(1 - D(\tilde{h}_i; \mathbf{X}, \mathbf{A}))$$

$h_i = \text{GNN}(X, A)$

$\tilde{h}_i$  negative sample

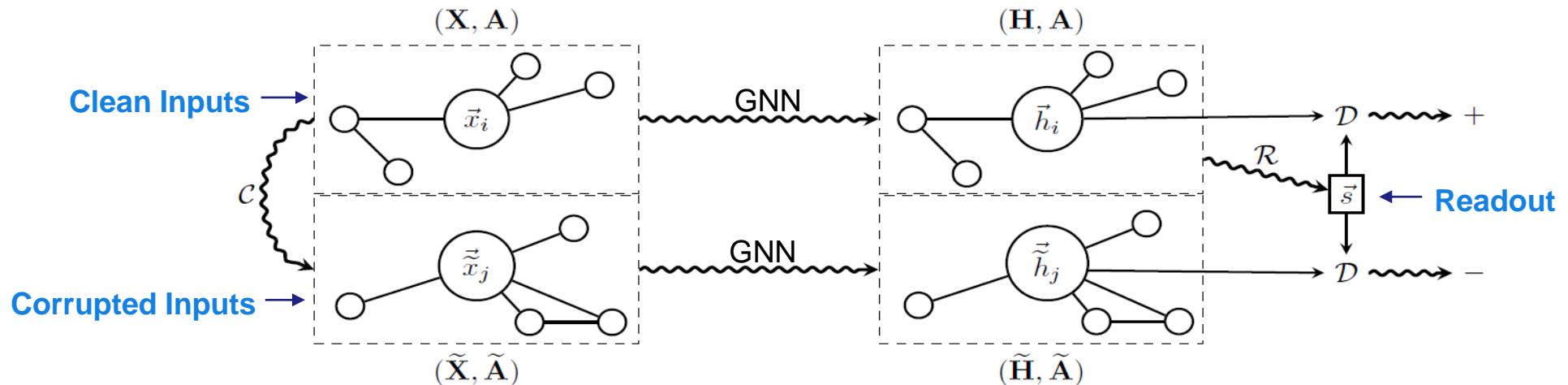
# Deep Graph Infomax (DGI)

It is hard to directly compute  $D(\tilde{h}_i; \mathbf{X}, \mathbf{A})$ , thus DGI resorts to readout  $\mathbf{s} = R(\mathbf{X}, \mathbf{A})$ :

$$\max_{\text{GNN}} I(X, A; h_i) \approx \max \log(D(h_i; \mathbf{X}, \mathbf{A})) + \log(1 - D(\tilde{h}_i; \mathbf{X}, \mathbf{A}))$$



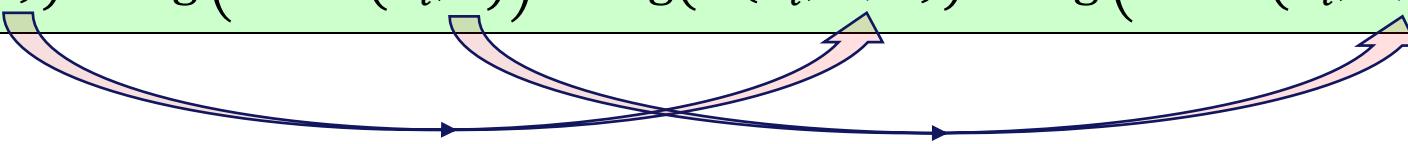
$$\max_{\text{GNN}} I(X, A; h_i) = \max_{\text{GNN}} \log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s}))$$



# Deep Graph Infomax (DGI)

It can be proved that, if the readout  $s = R(X, A)$  is injective,

$$\log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s})) = \log(D(h_i; \mathbf{X}, \mathbf{A})) + \log(1 - D(\tilde{h}_i; \mathbf{X}, \mathbf{A}))$$



It can be also proved that, if  $|\mathbf{X}| = |\mathbf{s}|$  is finite,

$$\max \log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s})) = \max I(h_i; X, A)$$



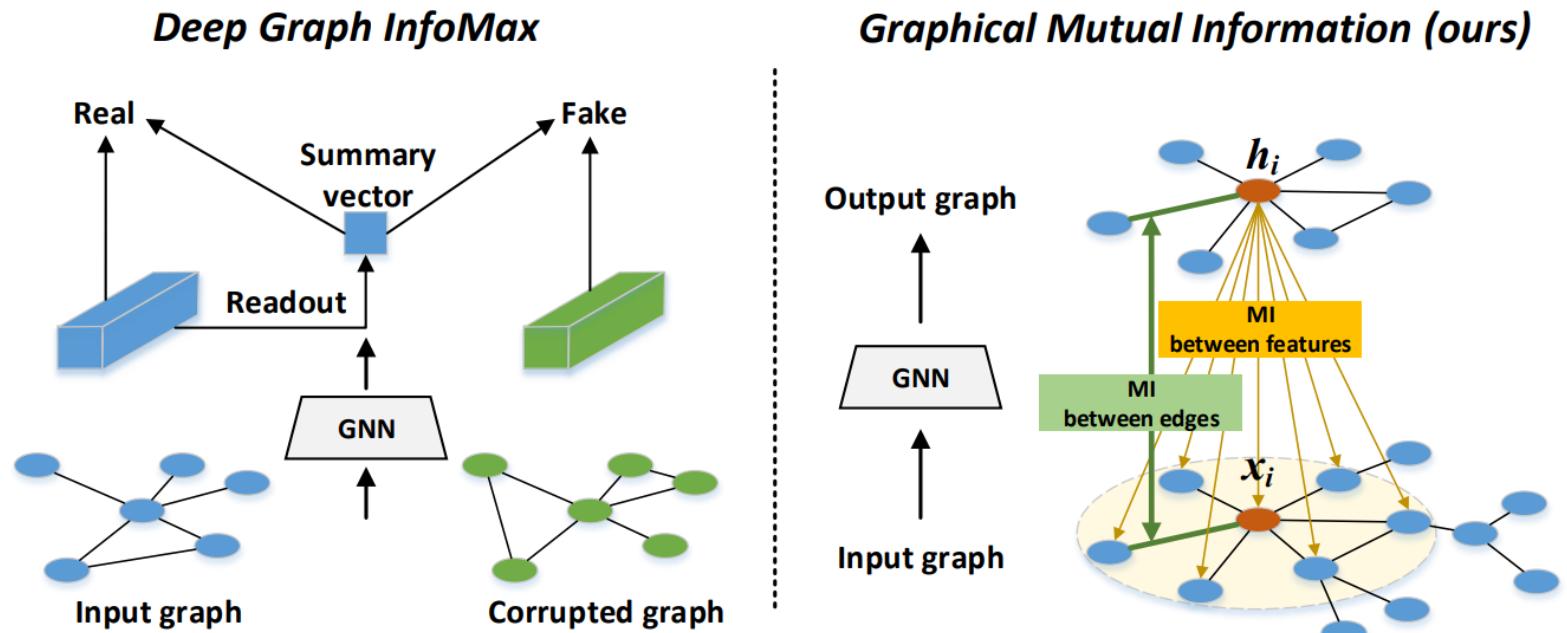
- Some issues in DGI

- Computing MI requires the injectivity of readout function
- It resorts to graph corruption to generate negative samples
- Distinct encoders and corruption functions for different tasks

Indirect  
Inefficient

# GMI: Graphical Mutual Information

- (Peng et al. 2020)



The basic idea of GMI is to compute the MI directly.

# GMI: Graphical Mutual Information

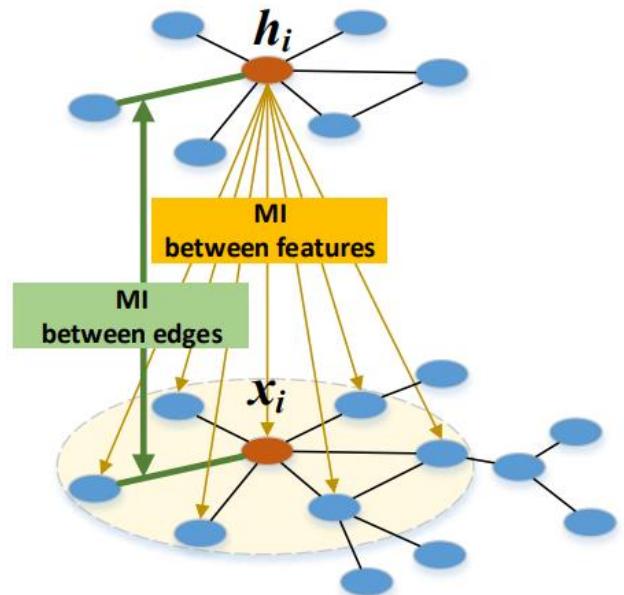
We define that,

$$I(X, A; h_i) \approx I(X; h_i) + \sum_{j \in N(i)} I(\sigma(h_i^T h_j); A_{ij})$$

Feature MI                      Topology MI

- It is both feature- and edge- aware;
- No need to readout or corruption;
- Feature MI can be further decomposed;

The basic idea of GMI is to compute the MI directly.





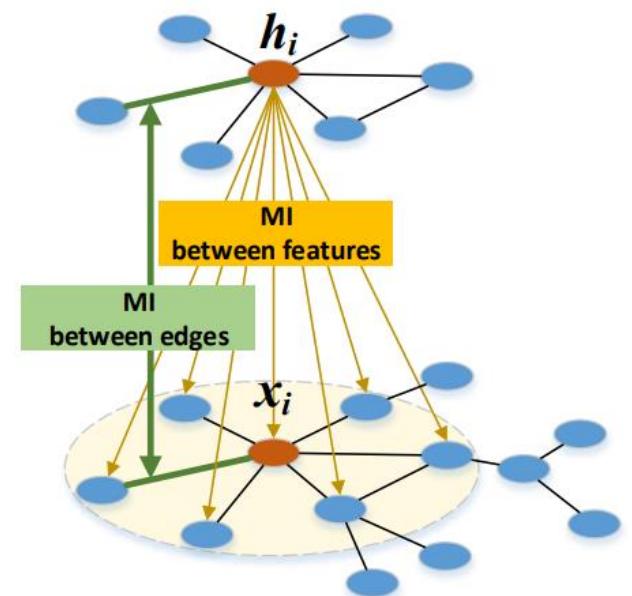
# GMI: Graphical Mutual Information

- (Peng et al. 2020)

It can be proved that, if certain mild condition meets,

$$I(X; h_i) = \sum_{j \in N(i)} w_{ij} I(x_j; h_i), \text{ for } 0 \leq w_{ij} \leq 1$$

The global MI is decomposed into a weighted sum of local MIs.  
It is not a bad idea to let  $w_{ij} = \sigma(h_i^T h_j)$



We then apply the JSD MI estimator to compute  $I(x_j; h_i)$  and  $I(\sigma(h_i^T h_j); A_{ij})$

# GMI: Graphical Mutual Information

- Node Classification

We use a universal backbone (GCN) for all tasks, different from DGI

Algorithm	Transductive			Inductive	
	Cora	Citeseer	PubMed	Reddit	PPI
EP-B loss	$79.4 \pm 0.1$	$69.3 \pm 0.2$	$78.6 \pm 0.2$	$93.8 \pm 0.03$	$61.8 \pm 0.04$
DGI loss	$82.2 \pm 0.2$	$72.2 \pm 0.2$	$78.9 \pm 0.3$	$94.3 \pm 0.02$	$62.3 \pm 0.02$
<b>FMI (ours)</b>	$78.3 \pm 0.1$	$72.0 \pm 0.2$	$79.1 \pm 0.3$	$94.7 \pm 0.03$	$64.8 \pm 0.03$
<b>GMI-mean (ours)</b>	<b><math>82.7 \pm 0.1</math></b>	<b><math>73.0 \pm 0.3</math></b>	<b><math>80.1 \pm 0.2</math></b>	<b><math>95.0 \pm 0.02</math></b>	<b><math>65.0 \pm 0.02</math></b>
<b>GMI-adaptive (ours)</b>	<b><math>83.0 \pm 0.3</math></b>	<b><math>72.4 \pm 0.1</math></b>	<b><math>79.9 \pm 0.2</math></b>	<b><math>94.9 \pm 0.02</math></b>	<b><math>64.6 \pm 0.03</math></b>

Codes: <https://github.com/zpeng27/GMI>

# GMI: Graphical Mutual Information

- Link Prediction

We use an universal backbone (GCN) for all tasks

<b>Algorithm</b>	<b>Cora</b>			<b>BlogCatalog</b>			<b>Flickr</b>			<b>PPI</b>
	<b>20.0%</b>	<b>50.0%</b>	<b>70.0%</b>	<b>20.0%</b>	<b>50.0%</b>	<b>70.0%</b>	<b>20.0%</b>	<b>50.0%</b>	<b>70.0%</b>	<b>22.7%</b>
DGI	95.6±0.3	94.6±0.4	94.4±0.2	77.2±0.4	76.4±0.4	75.5±0.3	90.3±0.3	89.0±0.4	74.1±0.7	77.4±0.1
<b>FMI (ours)</b>	<b>97.2±0.2</b>	<b>95.2±0.1</b>	<b>95.0±0.1</b>	<b>81.2±0.2</b>	<b>79.5±0.4</b>	<b>75.1±0.2</b>	<b>92.7±0.3</b>	<b>92.2±0.3</b>	<b>90.6±0.4</b>	<b>79.8±0.2</b>
<b>GMI (ours)</b>	<b>97.9±0.3</b>	<b>96.4±0.2</b>	<b>96.3±0.1</b>	<b>84.1±0.3</b>	<b>83.6±0.2</b>	<b>82.5±0.1</b>	<b>92.0±0.2</b>	<b>90.1±0.3</b>	<b>88.5±0.2</b>	<b>80.0±0.2</b>

Codes: <https://github.com/zpeng27/GMI>

# Summary

	Node-Classification	Link/Metapath Prediction	Graph-Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S <sup>2</sup> GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5] , GROVER [7] GCC [6]	
Information-based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]



Tencent  
AI Lab



# Break



5 min



Tencent  
AI Lab



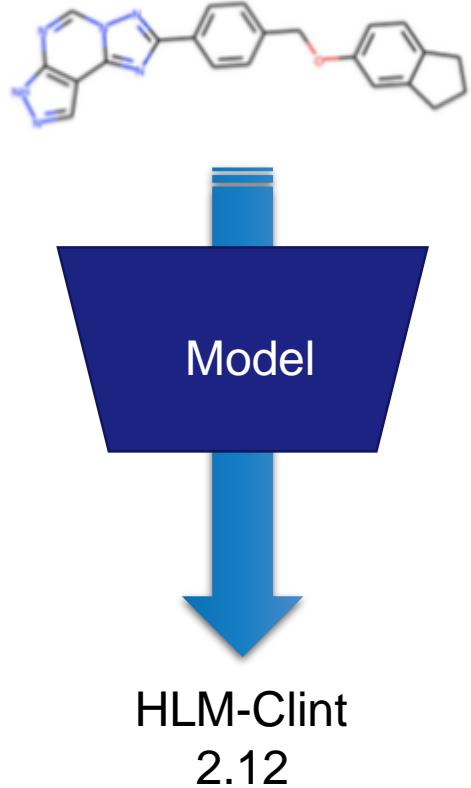
# Other Advanced Topics of GNNs

## Subgraph Recognition

# Problem of Subgraph Recognition

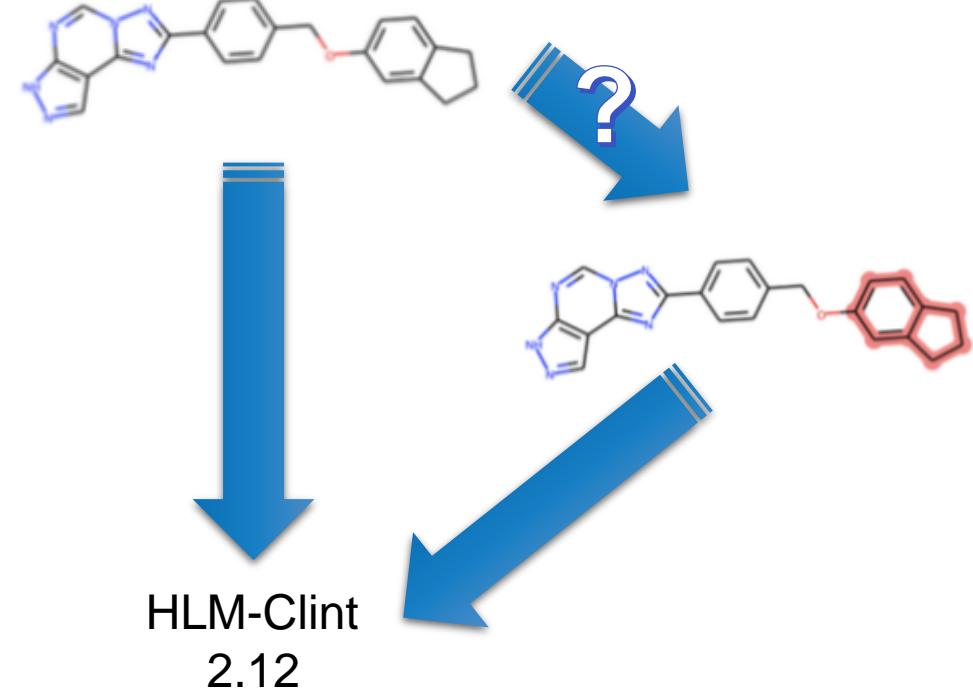
## Graph Classification

- Predict the property of a graph



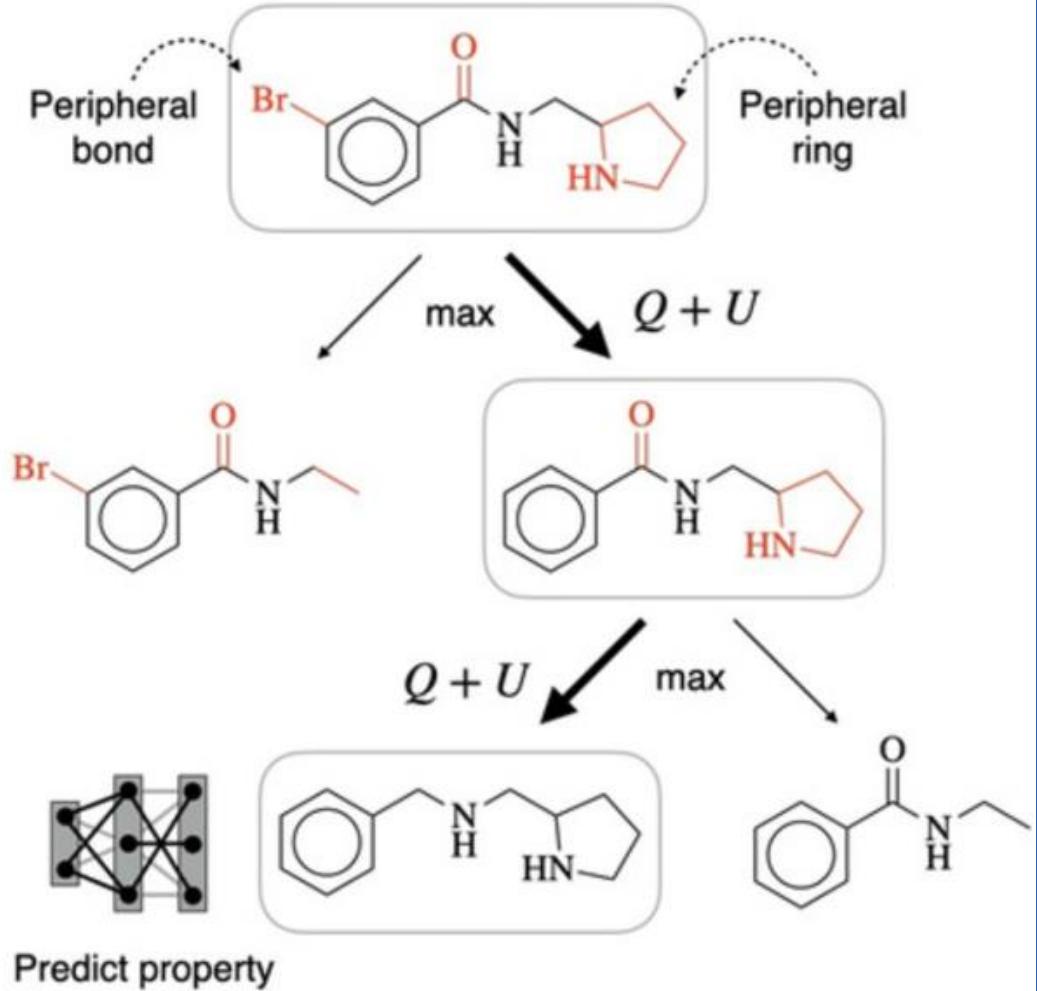
## Subgraph Recognition

- Recognize the sub-graph that most predict the property



# Monte Carlo Tree Search

MCTS (Jin et al. 2020)



## Solving Steps

- ◀ Iteratively remove peripheral bonds and rings to find subgraphs
- ◀ Evaluate each subgraph using the fixed property predictor
- ◀  $Q$  and  $U$  functions are MCTS parameters

## Weakness

- ◀ Need **many iterations** to delete components **step by step**
- ◀ Require **two stages** with a **fixed** property predictor
- ◀ Seek in a **large** search space

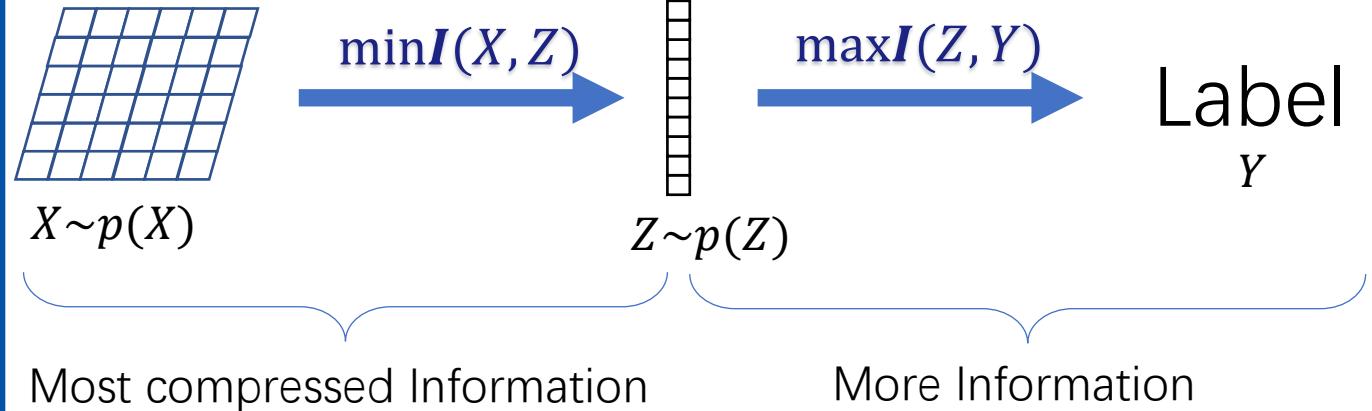
# Information Bottleneck

## Information Bottleneck

$$\max_Z I(Y, Z) - \beta I(X, Z)$$

- Where  $I(\cdot, \cdot)$  indicates the Mutual Information of two variables
- Compressed information between  $X$  and  $Z$
- Most informative  $Z$  to reflect the properties  $Y$

Traditional IB objective:  $\max_Z I(Y, Z) - \beta I(X, Z)$



# Variational Information Bottleneck

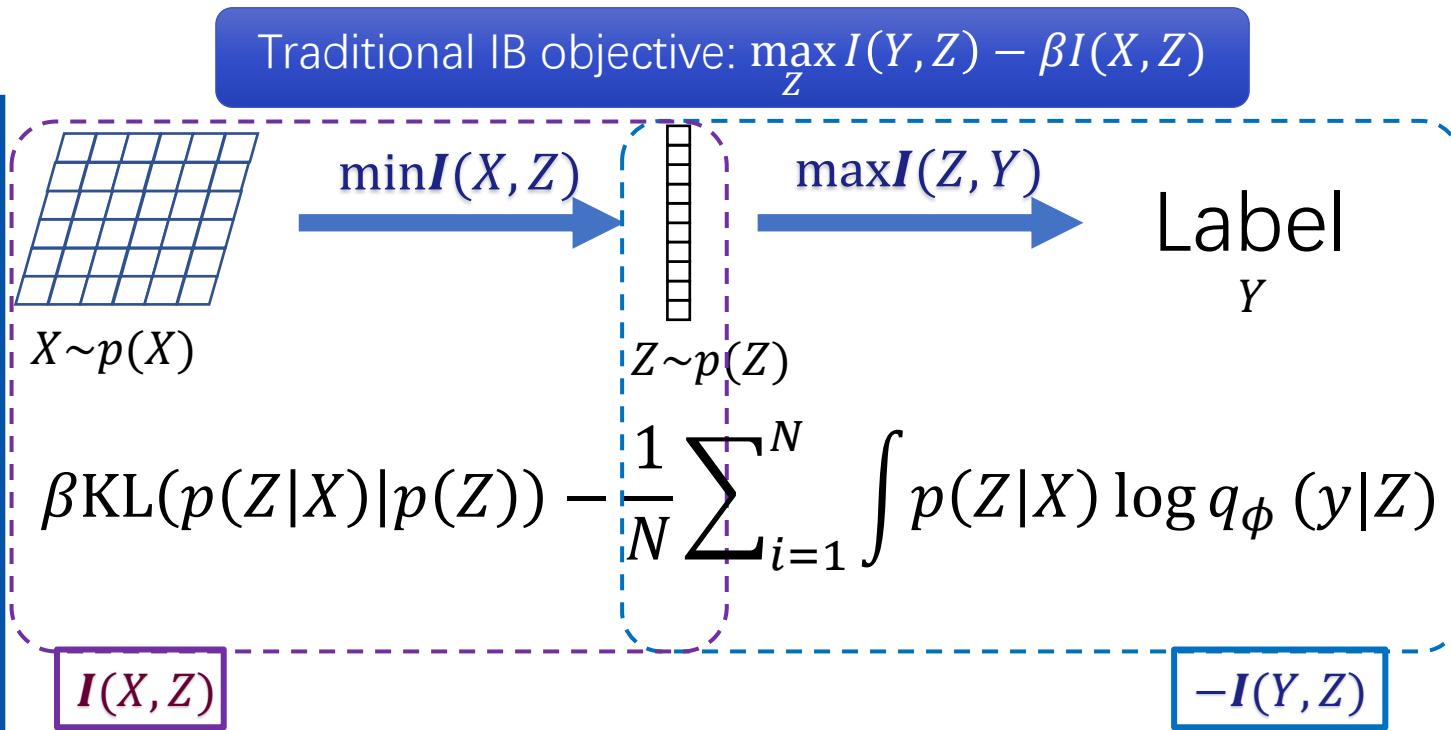
## Information Bottleneck

$$\max_Z I(Y, Z) - \beta I(X, Z)$$

- Where  $I(\cdot, \cdot)$  indicates the Mutual Information of two variables
- Solution: Variational information bottleneck

$$\mathcal{L}_{VIB} = \frac{1}{N} \sum_{i=1}^N \int p(Z|X) \log q_\phi(y|Z) - \beta \text{KL}(p(Z|X) \| p(Z))$$

Traditional IB objective:  $\max_Z I(Y, Z) - \beta I(X, Z)$



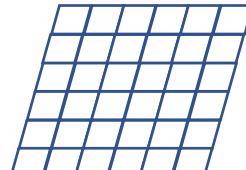
# Graph Information Bottleneck

Problem

Can we find an informative yet compressed subgraph,  $G_{sub}$ , from the graph,  $G$ ?

We leverage the idea from Information Bottleneck

Traditional IB objective:  $\max_z I(Y, Z) - \beta I(Z, X)$



$X \sim p(X)$

$\min I(X, Z)$

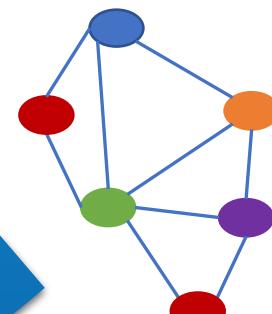


$Z \sim p(Z)$

$\max I(Z, Y)$

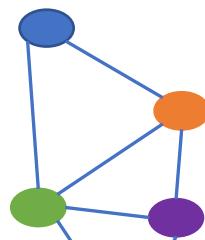
Label  
 $Y$

GIB objective:  $\max_{G_{sub}} I(Y, G_{sub}) - \beta I(G, G_{sub})$



$G \sim p(G)$

$\min I(G, G_{sub})$



$G_{sub} \sim p(G_{sub})$

$\max I(Y, G_{sub})$

Label  
 $Y$

Compression

Prediction

# Graph Information Bottleneck

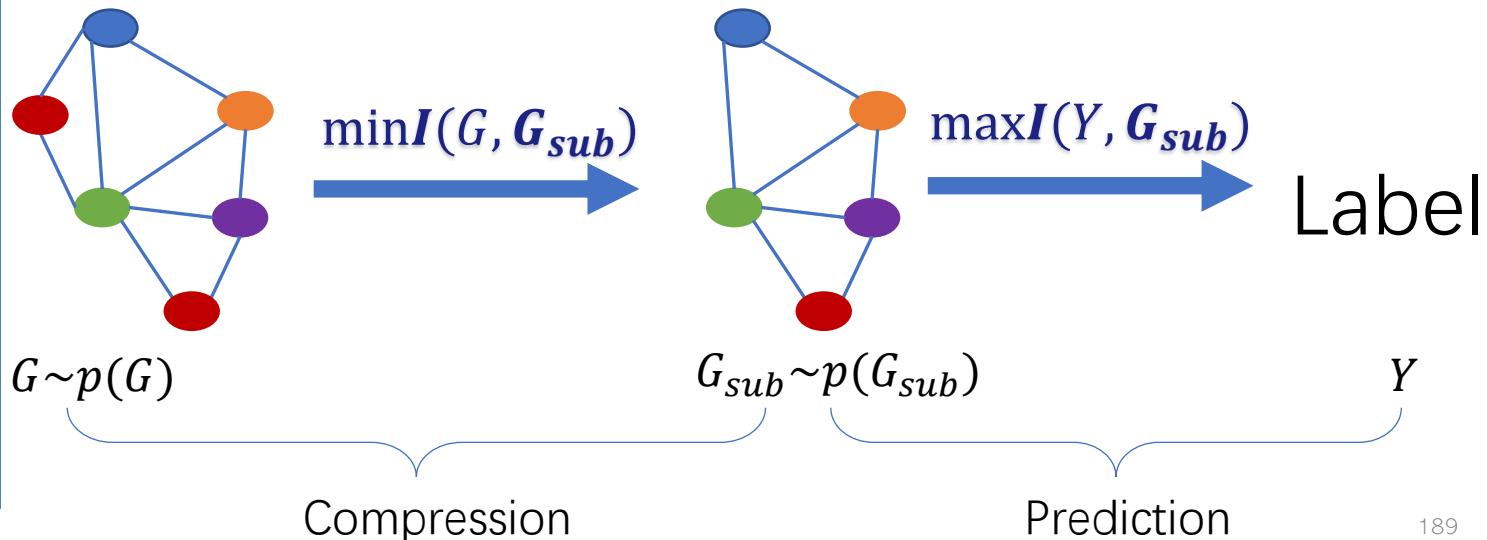
VIB to Graph?

- Use variational information bottleneck? **NO**

$$\frac{1}{N} \sum_{i=1}^N \int p(G_{sub}|G) \log q_\phi(y|G_{sub}) - \beta \text{KL}(p(G_{sub}|G) \| p(G_{sub}))$$

- The distribution of  $p(G_{sub})$  is unknown

GIB objective:  $\max_{G_{sub}} I(Y, G_{sub}) - \beta I(G, G_{sub})$

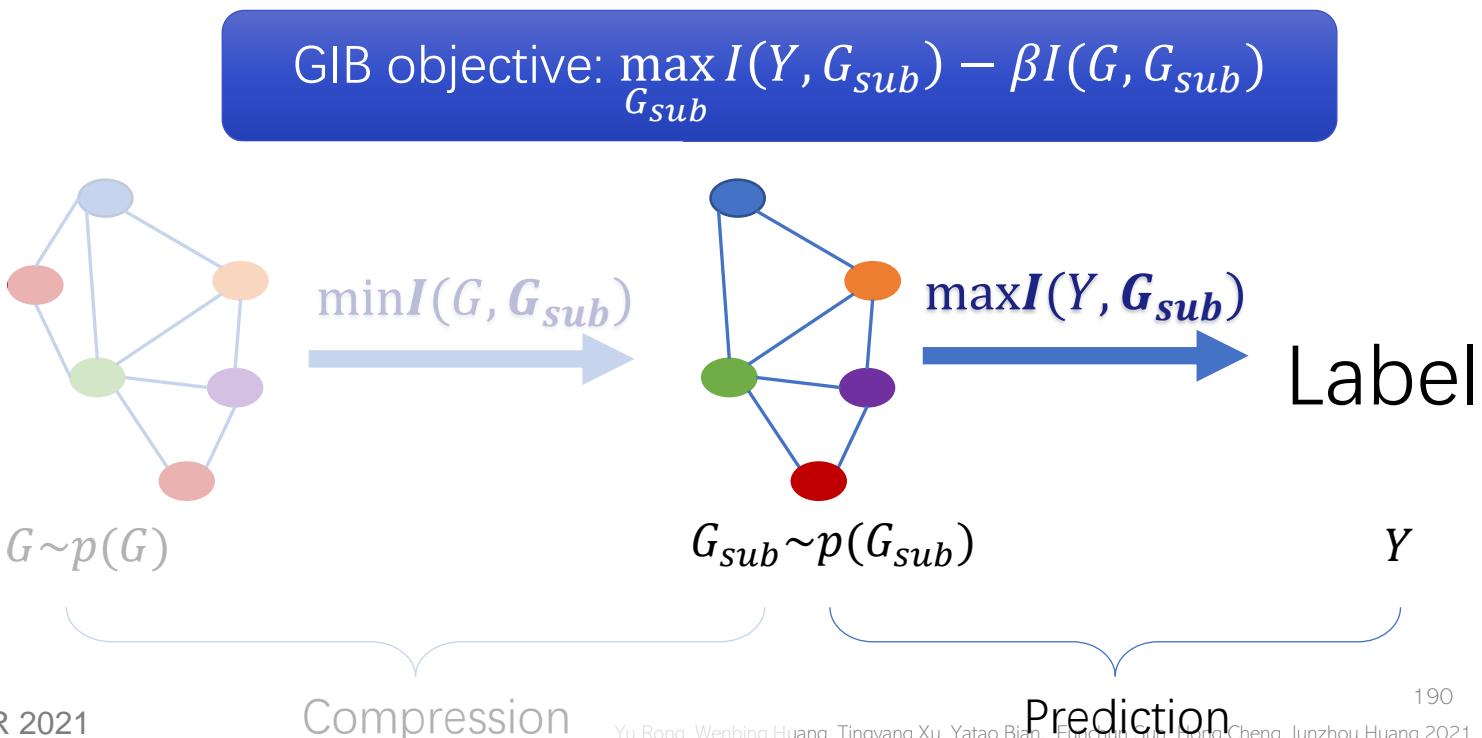


# GIB: Mutual Information between $Y$ and $G_{sub}$

$$I(Y, G_{sub})$$

- Maximize mutual information between the property and the sub-graph,  $I(Y, G_{sub})$

$$\begin{aligned} I(Y, G_{sub}) &\geq \frac{1}{N} \sum_{i=1}^N q_{\phi_1}(y|G_{sub}) \\ &=: -\mathcal{L}_{cls}(q_{\phi_1}(y|G_{sub}), y_{gt}) \end{aligned}$$



# GIB: Mutual Information between $G$ and $G_{sub}$

$$I(G, G_{sub})$$

- We fix  $G_{sub}$  to approximate  $\mathcal{L}_{MI}$  by  $f_{\phi_2}$
- DONSKER-VARADHAN representation of the KL-divergence (Sampling)

$$I(G, G_{sub}) = \max_{\phi_2} \mathcal{L}_{MI}(\phi_2, G_{sub})$$

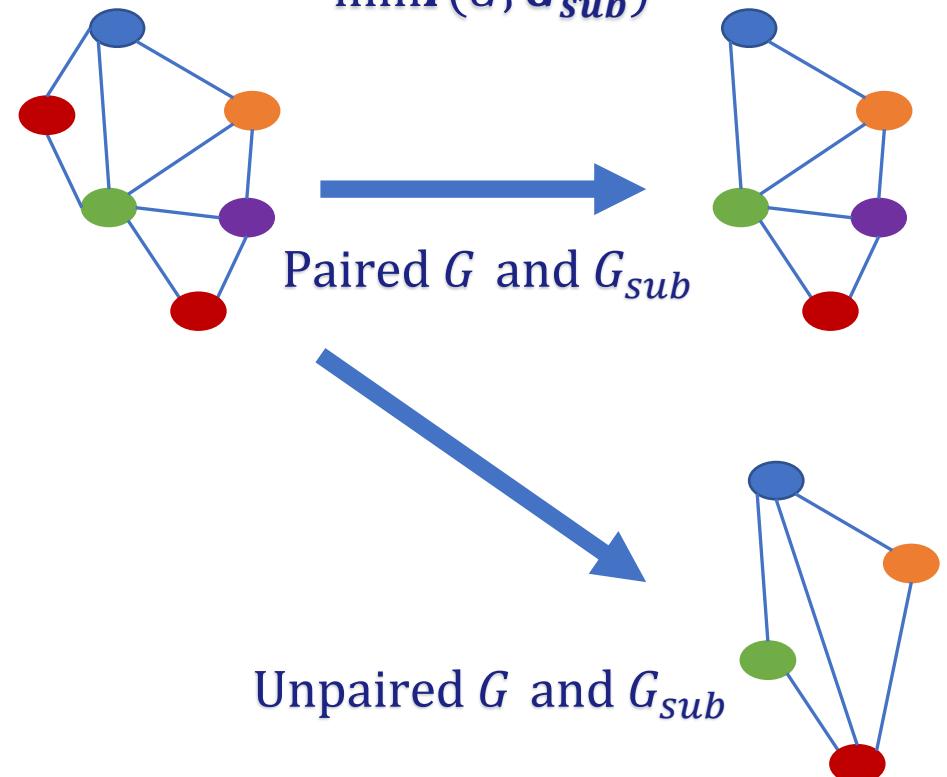
$$\approx \frac{1}{N} \sum_{i=1}^N f_{\phi_2}(G_i, G_{sub,i}) - \log \frac{1}{N} \sum_{i=1, j \neq i}^N e^{f_{\phi_2}(G_i, G_{sub,j})}$$

Paired  $G$  and  $G_{sub}$  in a batch

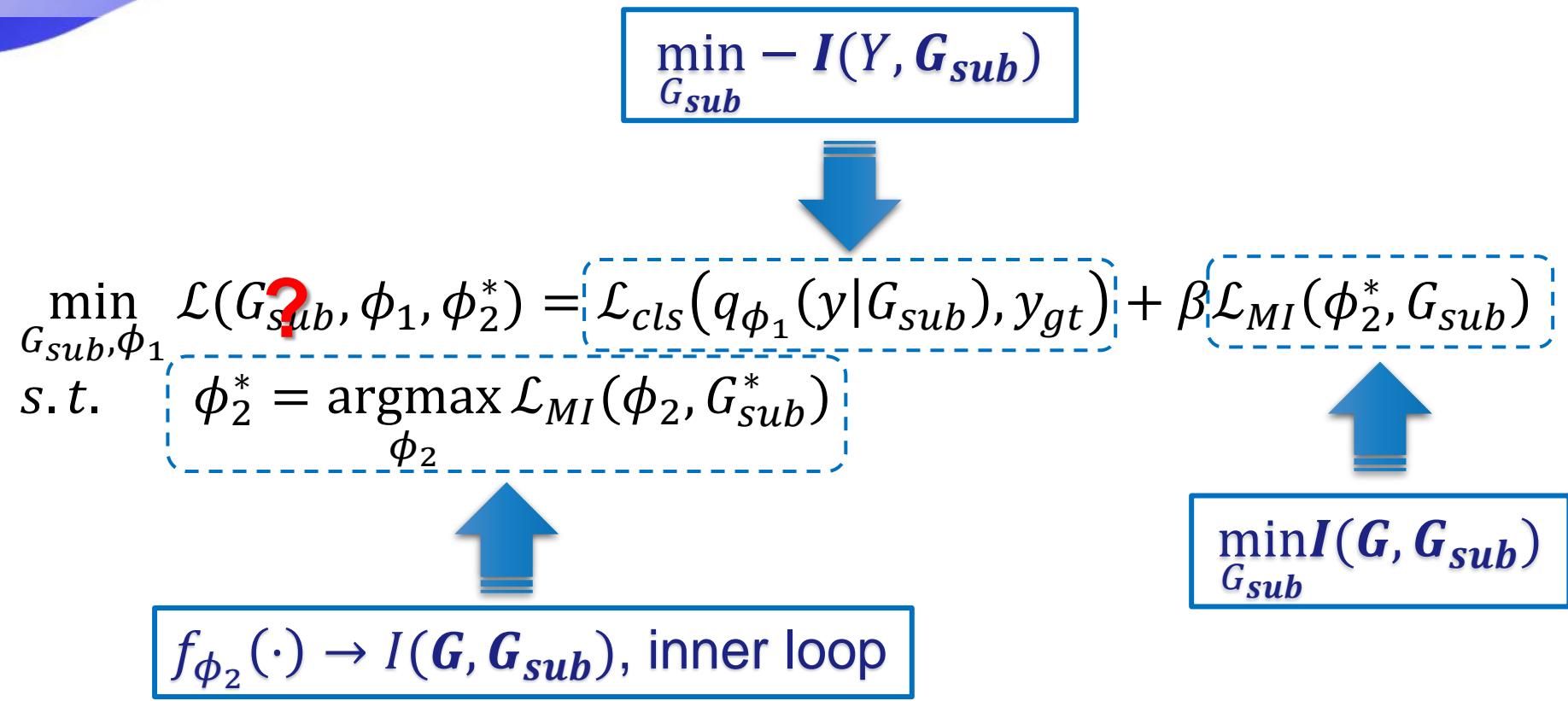
Unpaired  $G$  and  $G_{sub}$  in a batch

GIB objective:  $\max_{G_{sub}} I(Y, G_{sub}) - \beta I(G, G_{sub})$

$$\min I(G, G_{sub})$$



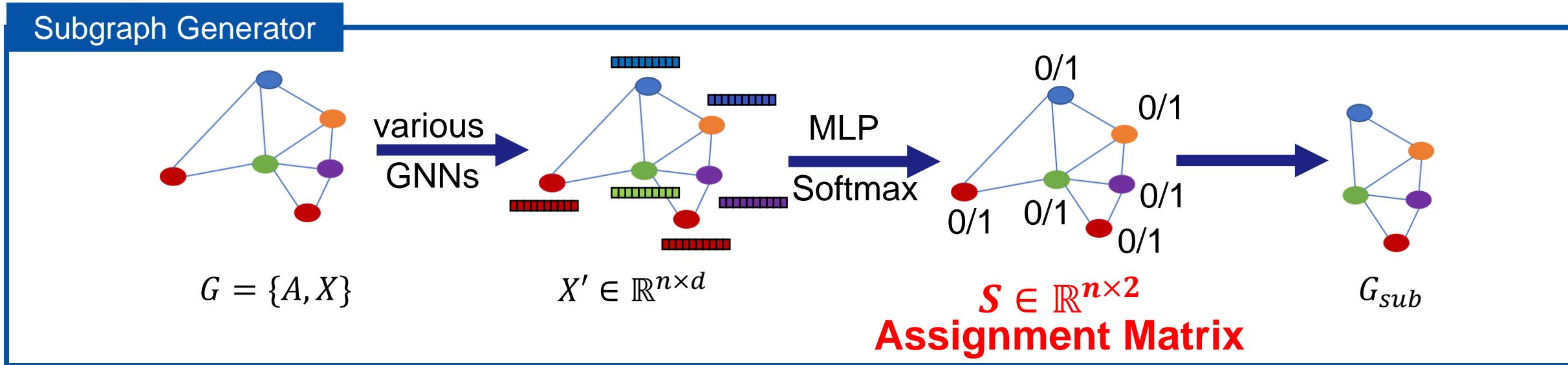
# GIB: Overall Objective Function



Problem

How to dynamically learn a  $G_{sub}$  for the graph,  $G$ ?

# GIB: How to get $G_{sub}$



## Connectivity Loss

$$\mathcal{L}_{con} = ||\text{Norm}(S^T AS) - I_2||_F$$

where *Norm* denotes row-wise normalization. This loss regularize:

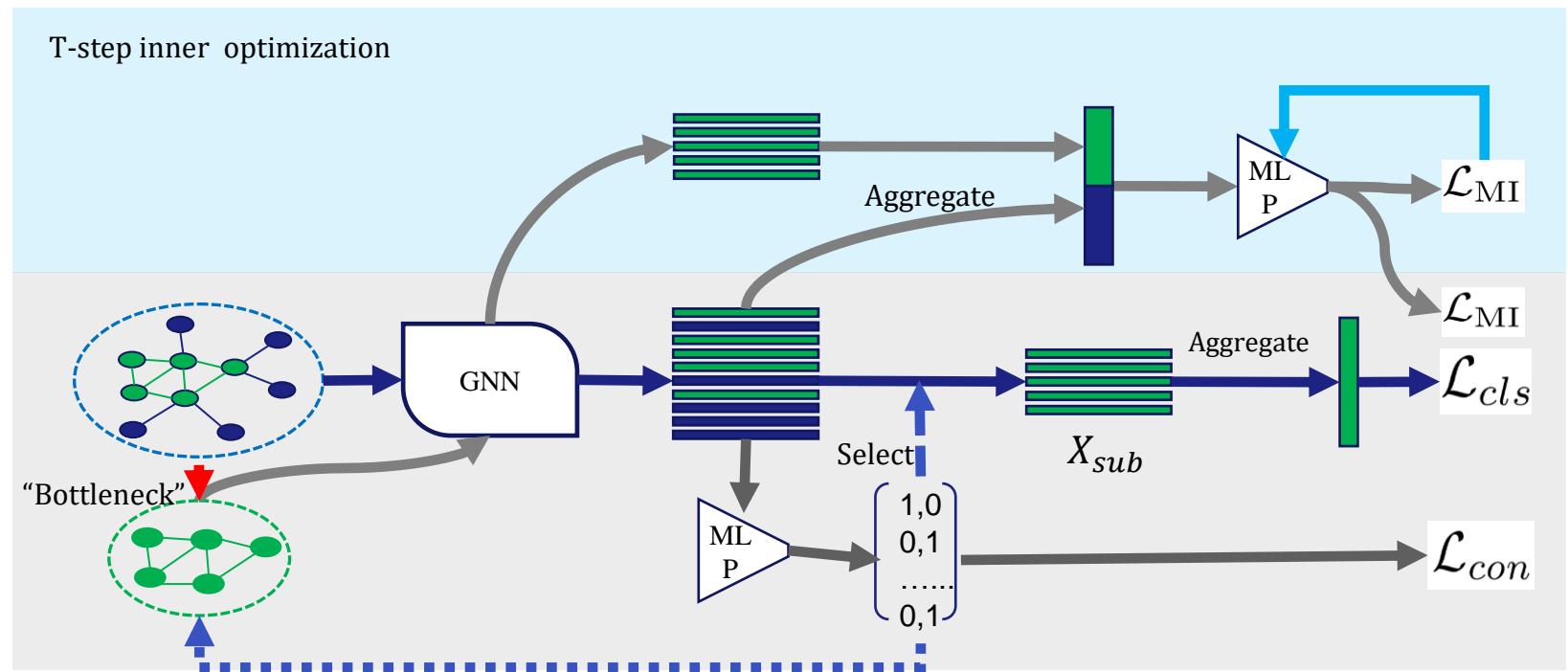
- ◀ Select **distinctive nodes**;
- ◀ Get **similar assignment** in neighborhoods;
- ◀ Avoid **trivial solutions** that all nodes are assigned to  $G_{sub}$  or  $\overline{G_{sub}}$ .

# Graph Information Bottleneck

## Objective Function

$$\begin{aligned} \min_{G_S, \phi_1} \mathcal{L}(G_S, \phi_1, \phi_2^*) &= \mathcal{L}_{cls}(q_{\phi_1}(y|G_S), y_{gt}) + \beta \mathcal{L}_{MI}(\phi_2^*, G_S) + \mathcal{L}_{con}(S) \\ s.t. \quad \phi_2^* &= \operatorname{argmax}_{\phi_2} \mathcal{L}_{MI}(\phi_2, G_S^*) \end{aligned}$$

## Flowchart

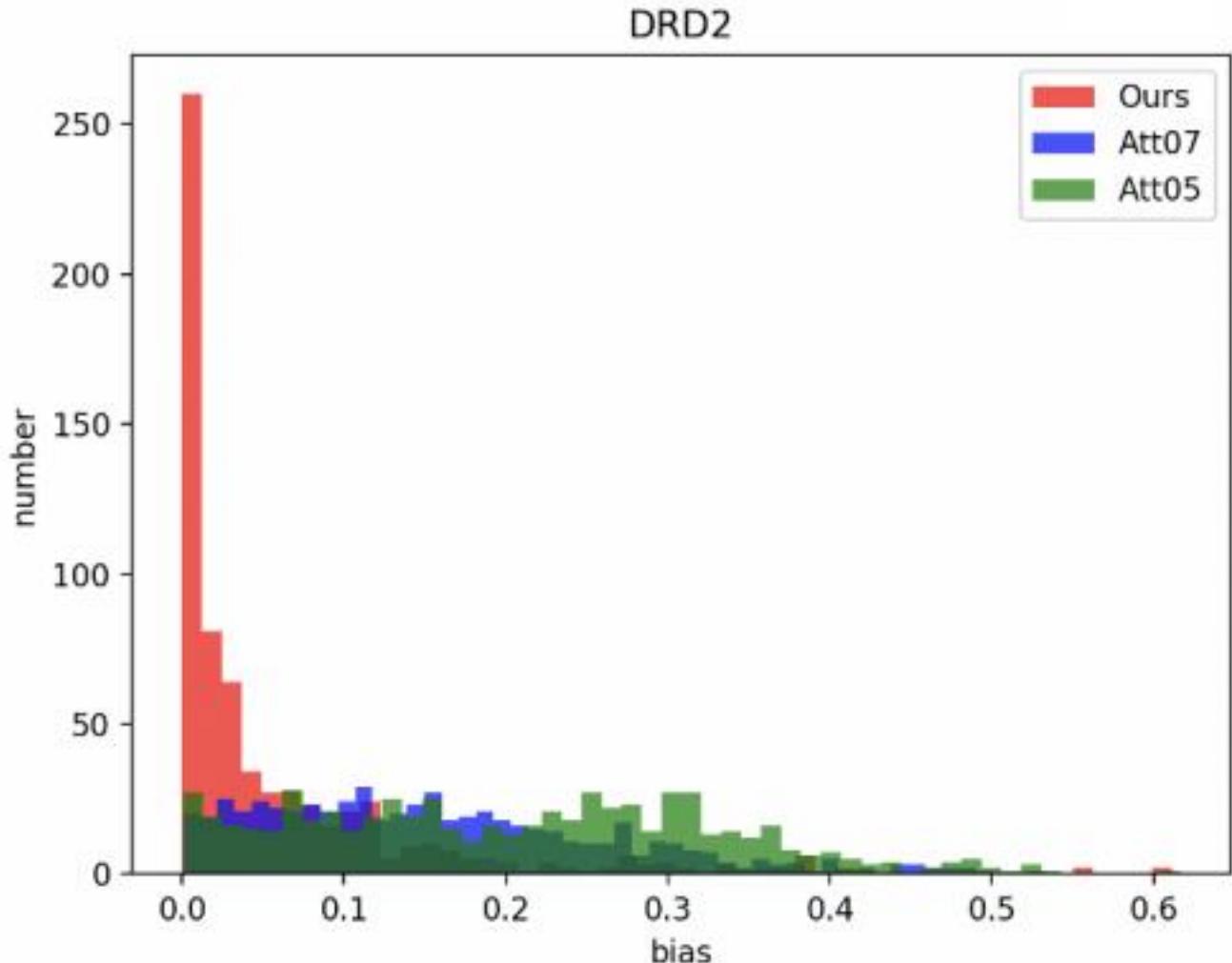


# Improve Graph Classification

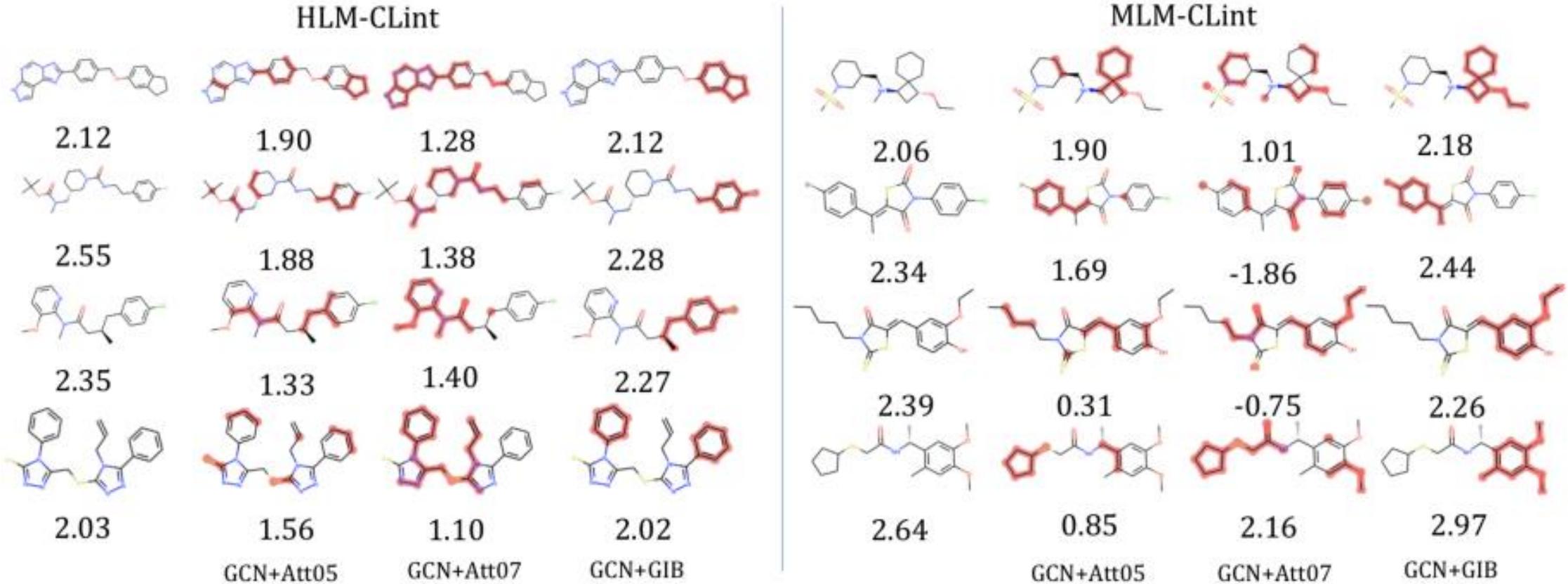
Method	MUTAG	PROTEINS	IMDB-BINARY	DD
SortPool	<b>0.844 ± 0.141</b>	0.747 ± 0.044	0.712 ± 0.047	0.732 ± 0.087
ASAPool	0.743 ± 0.077	0.721 ± 0.043	0.715 ± 0.044	0.717 ± 0.037
DiffPool	0.839 ± 0.097	0.727 ± 0.046	0.709 ± 0.053	0.778 ± 0.030
EdgePool	0.759 ± 0.077	0.723 ± 0.044	0.728 ± 0.044	0.736 ± 0.040
AttPool	0.721 ± 0.086	0.728 ± 0.041	0.722 ± 0.047	0.711 ± 0.055
GCN	0.743±0.110	0.719±0.041	0.707 ± 0.037	0.725 ± 0.046
GraphSAGE	0.743±0.077	0.721 ± 0.042	0.709 ± 0.041	0.729 ± 0.041
GIN	0.825±0.068	0.707 ± 0.056	0.732 ± 0.048	0.730 ± 0.033
GAT	0.738 ± 0.074	0.714 ± 0.040	0.713 ± 0.042	0.695 ± 0.045
GAT + DropEdge	0.743±0.081	0.711±0.043	0.710±0.041	0.717±0.035
<b>GCN+GIB</b>	0.776 ± 0.075	0.748 ± 0.046	0.722 ± 0.039	0.765 ± 0.050
<b>GraphSAGE+GIB</b>	0.760 ± 0.074	0.734 ± 0.043	0.719 ± 0.052	<b>0.781 ± 0.042</b>
<b>GIN+GIB</b>	0.839 ± 0.064	<b>0.749 ± 0.051</b>	<b>0.737 ± 0.070</b>	0.747 ± 0.039
<b>GAT+GIB</b>	0.749 ± 0.097	0.737 ± 0.044	0.729 ± 0.046	0.769 ± 0.040
<b>GAT+GIB+DropEdge</b>	0.754±0.085	0.737±0.037	0.731±0.003	0.776±0.034

# Consistent Prediction between Graph and Subgraph

$$\Delta_{pred} = |Pred(\text{Chemical Graph}) - Pred(\text{Subgraph})|$$



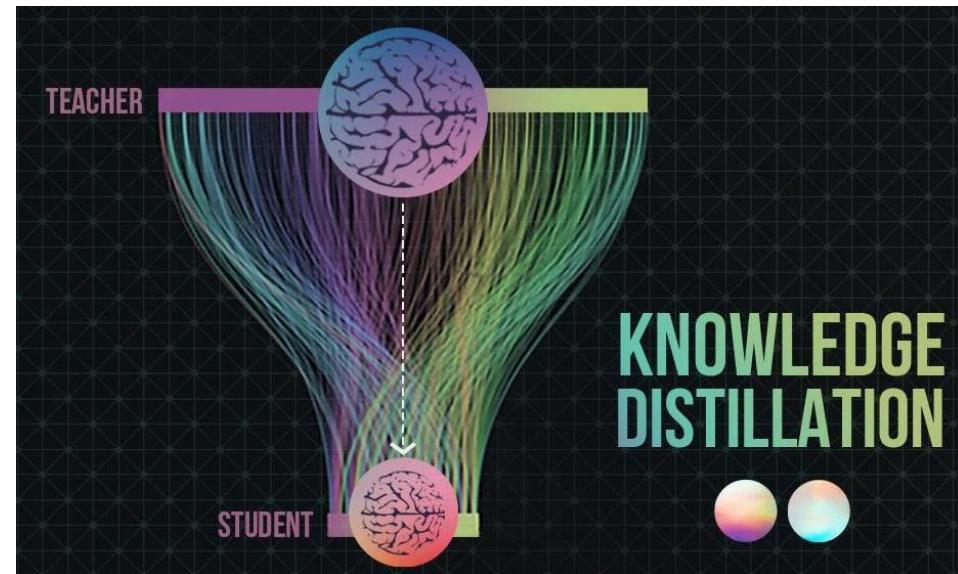
# Interpretable Subgraphs



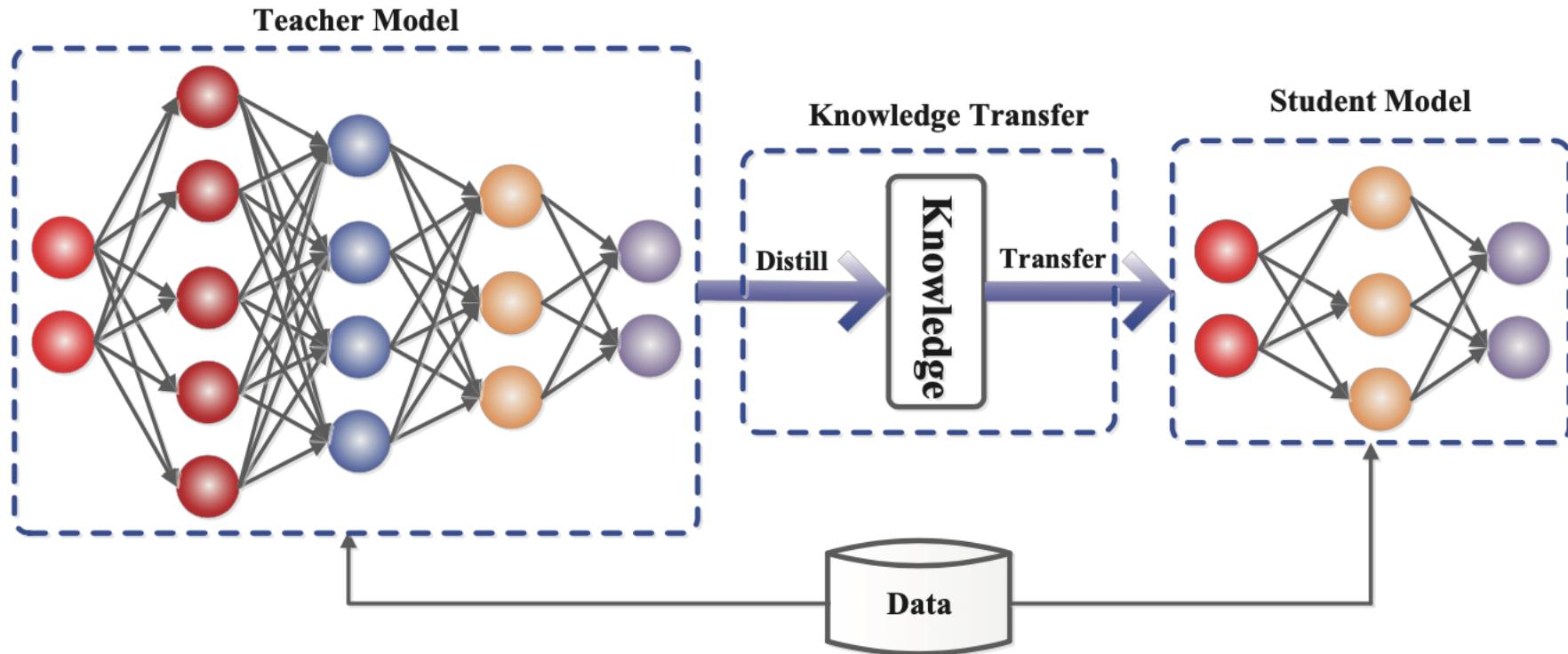
HLM-CLint & MLM-CLint: vitro human and mouse liver microsome metabolic stability (base 10 logarithm of mL/min/g)

# Other Advanced Topics of GNNs

Knowledge Distillation for GNNs



# Knowledge Distillation on CNNs



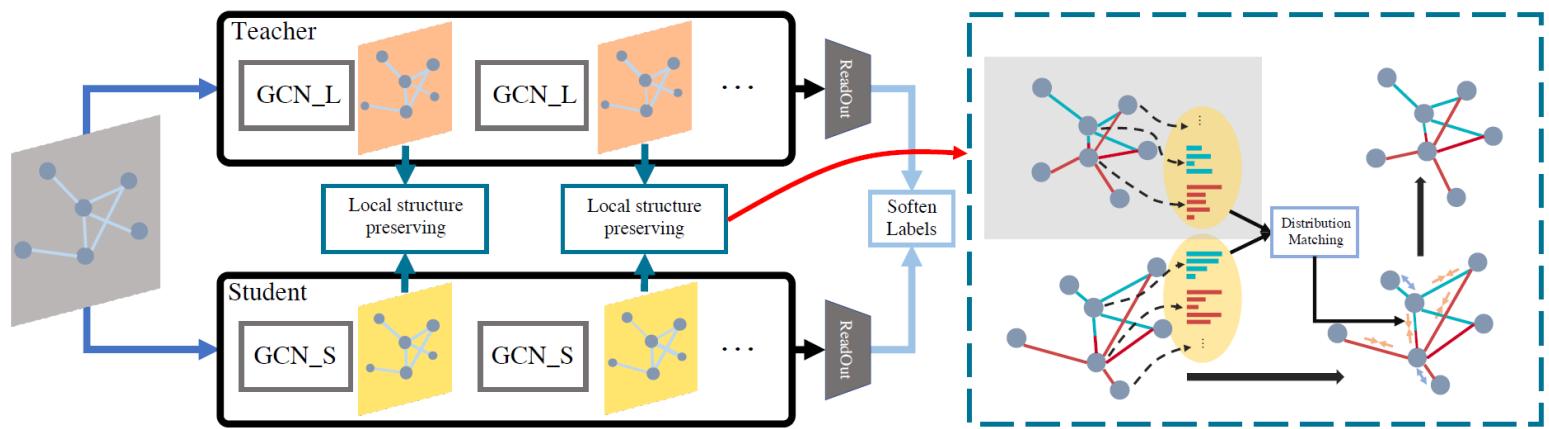
Knowledge Distillation (KD) has demonstrated its effectiveness on improving compact and lightweight Convolutional Neural Networks (CNNs).

Hinton, Geoffrey et al. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015)  
Gou, Jianping, et al. "Knowledge distillation: A survey." International Journal of Computer Vision (2021): 1-31.



# Teacher-Student Knowledge Distillation on GNNs

Recently, Yang et al. proposed to combine KD with the training of Graph Neural Networks (GNNs).



However, the conventional teacher-student knowledge distillation framework leads to two major concerns:

1. Training inefficiency
2. Sophisticated design of a qualified teacher model

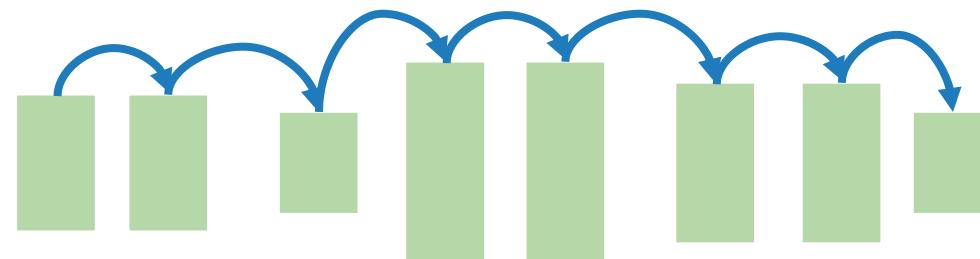
# Why Self-distilling Graph Neural Networks?

## 💡 Knowledge Distillation on GNNs → Self-Distillation on a single GNN

- |   |   |   |
|---|---|---|
| (1) Training Inefficiency;<br>(2) Difficulty on searching for a qualified teacher |  | (1) Marginal cost increase;<br>(2) Perform knowledge extraction and transfer among layers |
|---|---|---|

## 💡 How to conduct self KD in GNNs?

- 💡 Nodes tend to lose information when GNNs go deeper
- 💡 can distill some knowledge from **shallow** GNN layers!





# Knowledge Definition: Neighborhood Discrepancy Rate (NDR)

**Naïve** solution: define edge-wise difference:  $\|\mathbf{x}_{v^*}^{(l)} - \mathbf{x}_c^{(l)}\|_p$

## Proposition

$$\forall \epsilon > 0, \left\| \mathbf{x}_{v^*}^{(l)} - (\mathbf{D}^{-1} \mathbf{A} \mathbf{X})_{v^*}^{(l)} \right\|_p < \sum_{c \in N(v)} \left\| \mathbf{x}_{v^*}^{(l)} - \mathbf{x}_c^{(l)} \right\|_p < \epsilon, \text{ when } l \rightarrow \infty$$

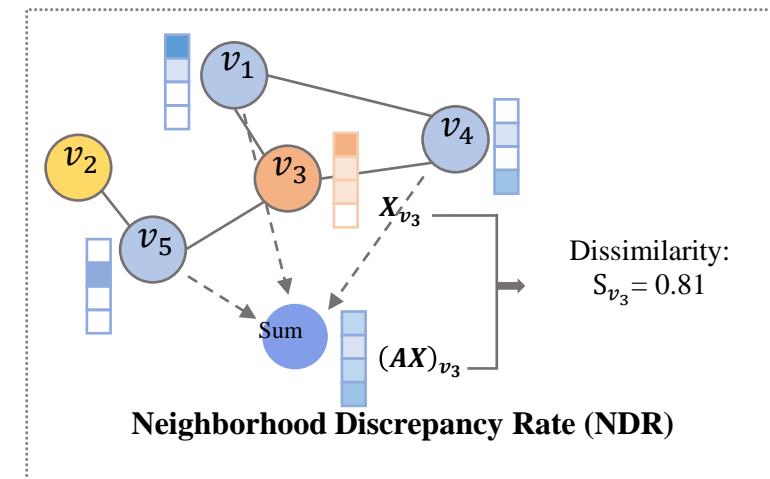
## Knowledge Definition

- Neighborhood Discrepancy Rate (NDR) of node  $v$ :

$$S_v^{(l)} = 1 - \frac{\mathbf{x}_{v^*}^{(l)} (\mathbf{A} \mathbf{X}^{(l)})_{v^*}^T}{\|\mathbf{x}_{v^*}^{(l)}\|_2 \|(\mathbf{A} \mathbf{X}^{(l)})_{v^*}\|_2}$$

- NDR vector of the graph:

$$\mathbf{s}^{(l)} = (S_1^{(l)}, \dots, S_N^{(l)})$$

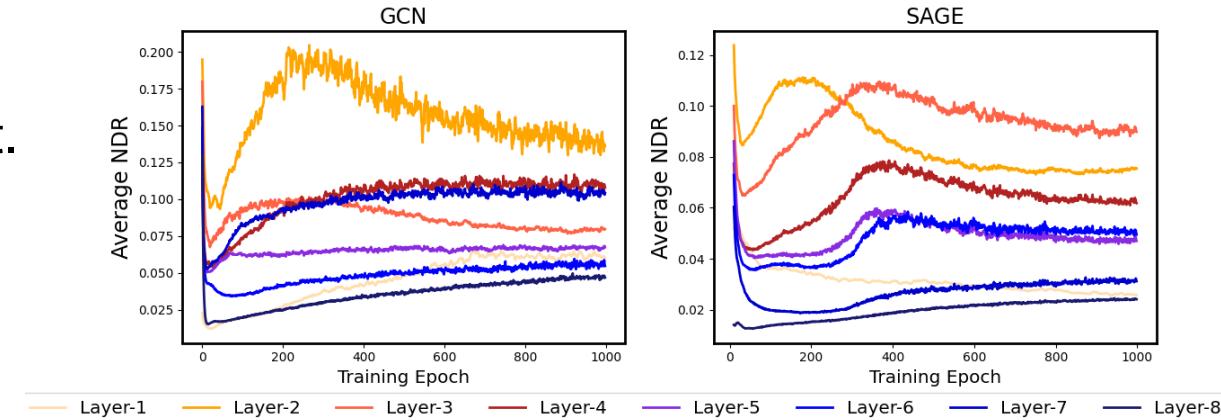


# Adaptive Distilling Strategy: Adaptive Discrepancy Retaining (ADR)

## Adaptive Discrepancy Retaining (ADR)

- ▶ Progressive (layer-by-layer) transfer.
- ▶ Adaptively select the initial supervision target.
  - ▶  $l^* = \operatorname{argmax}_{k \in [L-1]} \|\mathbf{s}^{(k)}\|$
- ▶ Explicit teacher selection.
- ▶ Regularizer formulation:

$$\mathcal{L}_N = \sum_{l=l^*, \dots, L-1} 1(\|\mathbf{s}^{(l)}\| > \|\mathbf{s}^{(l+1)}\|) d^2(\mathbf{s}^{(l+1)}, \mathbf{s}^{(l)}), \quad d^2(\mathbf{s}^{(l+1)}, \mathbf{s}^{(l)}) = \|\mathbf{s}^{(l+1)} - \text{StopGradient}(\mathbf{s}^{(l)})\|_2^2$$

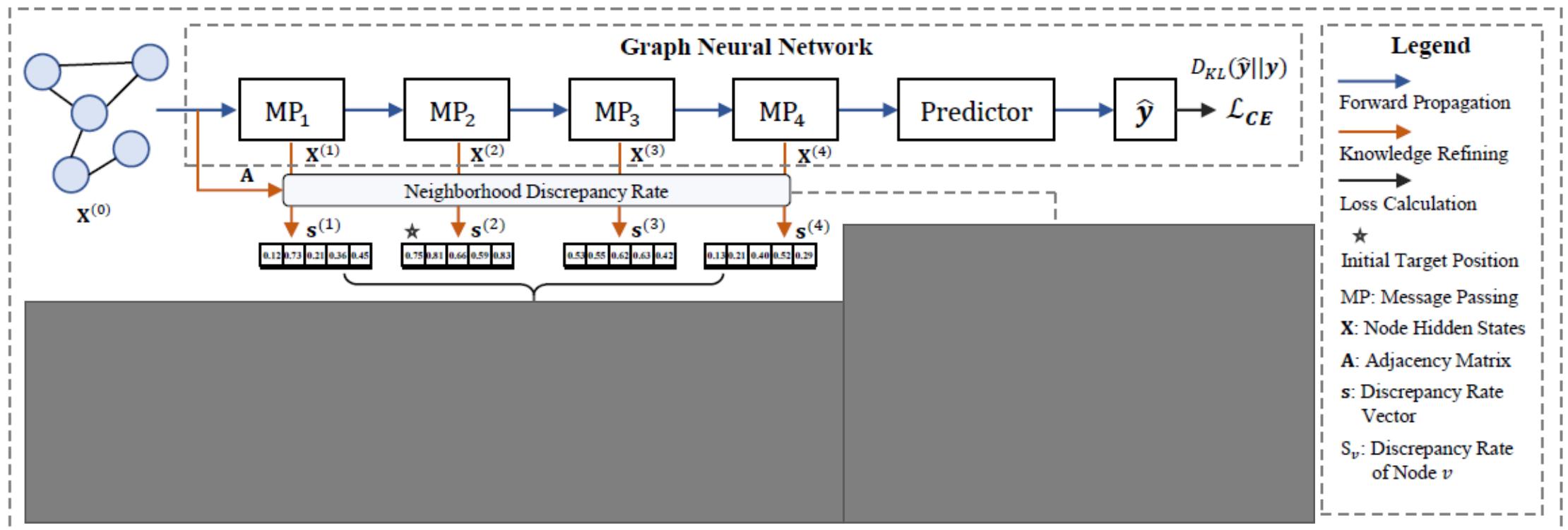


- ▶ Generic GNN Self-Distillation Framework by considering various granularities
  - ▶ knowledge over arbitrary subgraph
  - ▶ knowledge over the entire graph

# Overall GNN Self-KD Framework

## 💡 (Teacher-free) Self-Distillation on GNNs

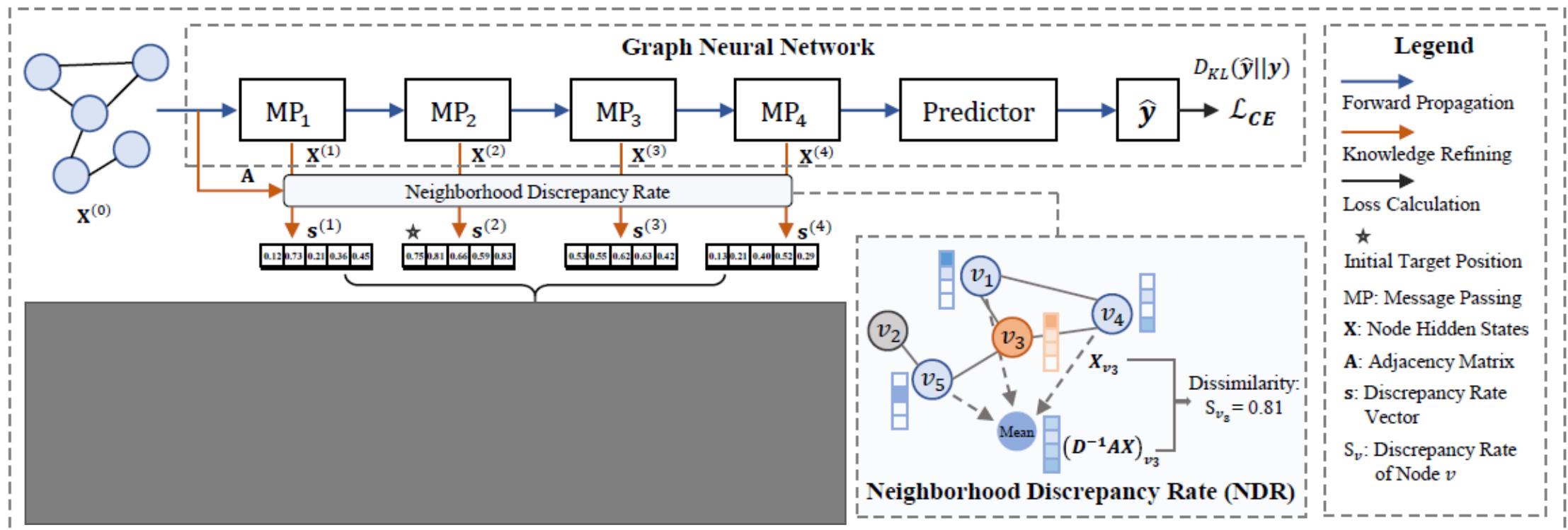
- 💡 (1) Marginal cost increase; (2) Perform knowledge extraction and transfer among layers



# Overall GNN Self-KD Framework

## 💡 (Teacher-free) Self-Distillation on GNNs

- 💡 (1) Marginal cost increase; (2) Perform knowledge extraction and transfer among layers

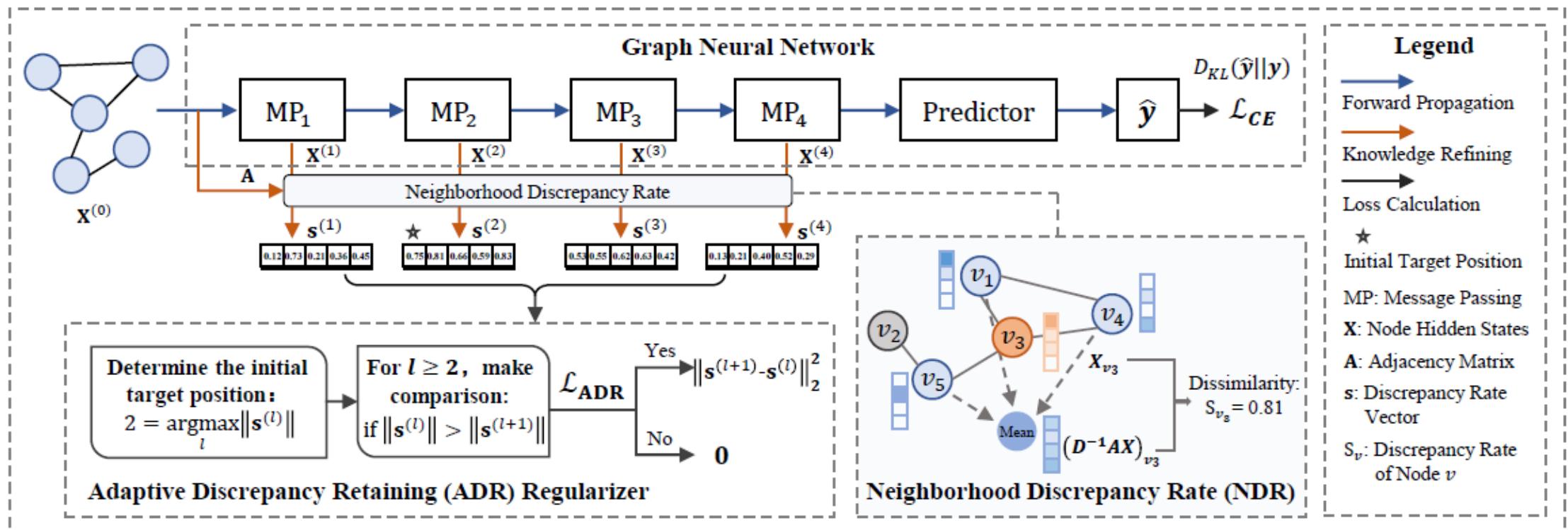




# Overall GNN Self-KD Framework

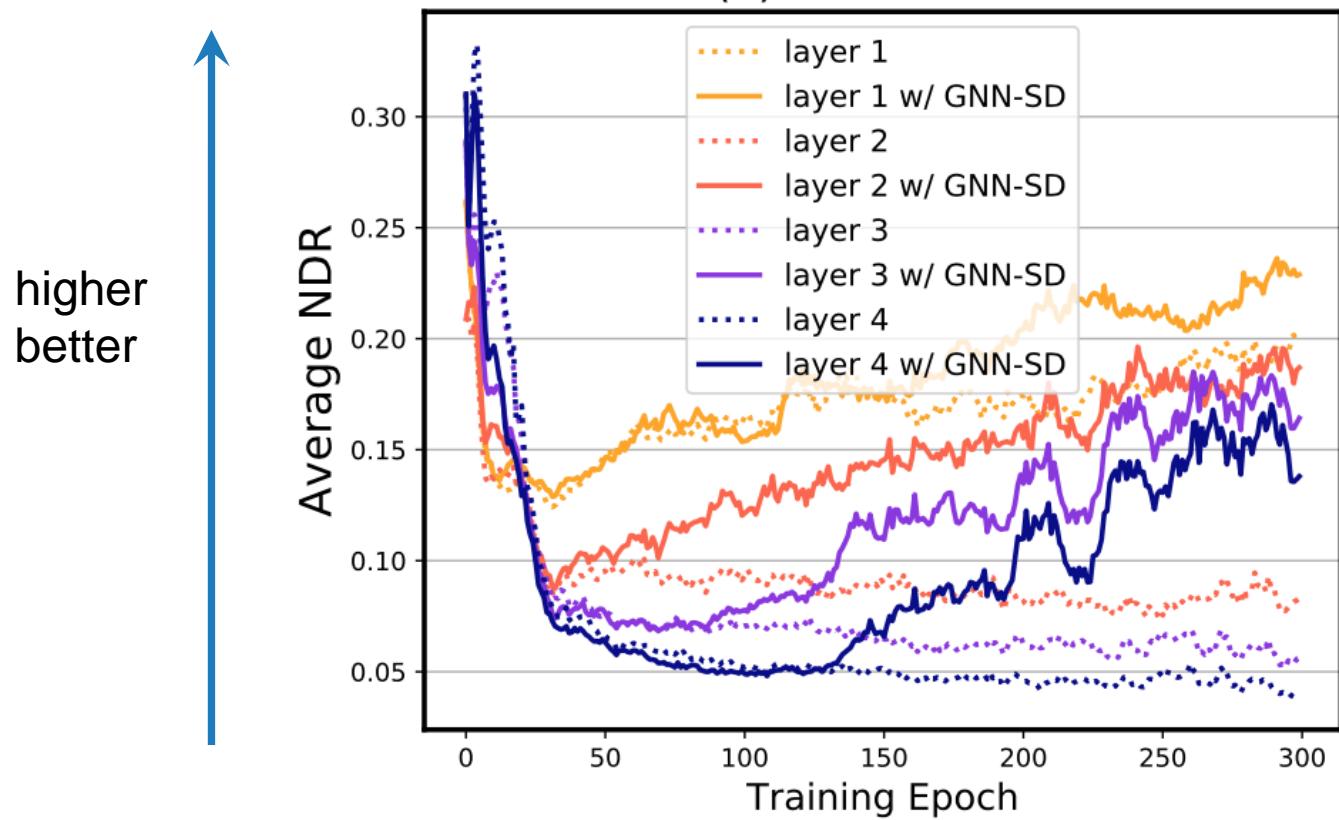
## 💡 (Teacher-free) Self-Distillation on GNNs

- 💡 (1) Marginal cost increase; (2) Perform knowledge extraction and transfer among layers



# The Effect of the NDR Regularizer

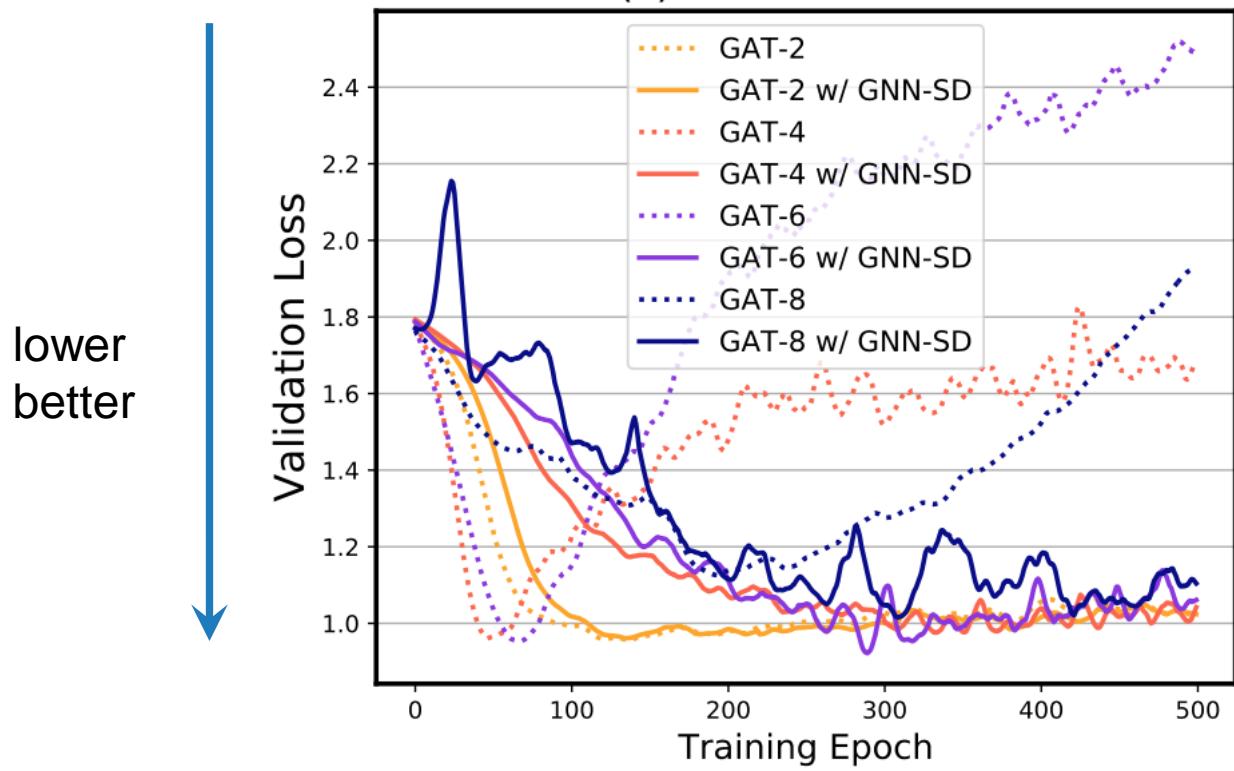
NDR between training w/ (solid line) and w/o (dotted) the GNN-SD strategy



The self-distillation strategy could improve the NDR for different layers

# GNN-SD Improves Generalization Performance

GAT-n: GAT with n layers  
dotted lines: w/o GNN-SD  
solid lines: w GNN-SD



The self-distillation strategy constantly improves the generalization performance by achieving lower validation losses



Tencent  
AI Lab



# Applications

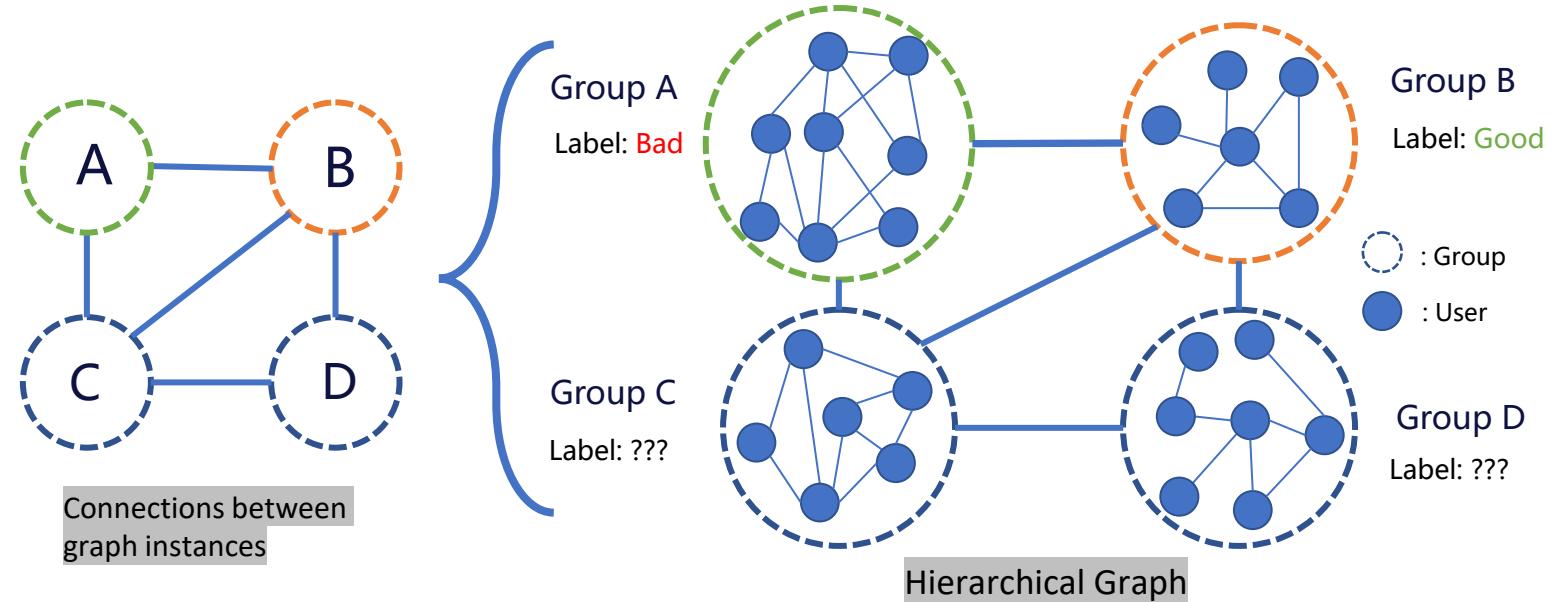
## GNN in Social Networks

# GNN in Social Networks

- "Semi-supervised graph classification: A hierarchical graph perspective." **WWW 2019**
- "Inverse Graph Identification: Can We Identify Node Labels Given Graph Labels?" **arXiv 2020**

# Hierarchical Graph Classification

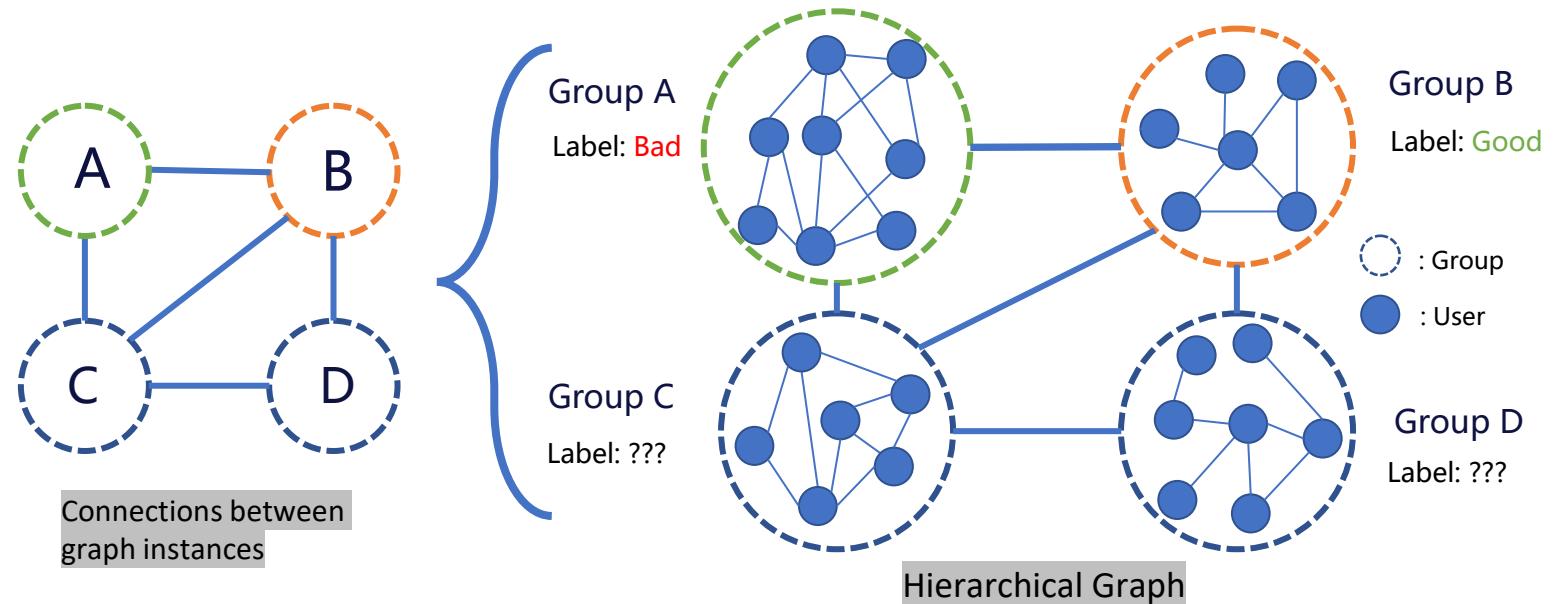
- **Hierarchical Graph:** A set of graph instances are interconnected via edges.
  - Social network with group structure.
  - Document (graph-of-words) collection with citation relation.



- **The Problem:** predicts the class label of graph instances in a hierarchical graph.

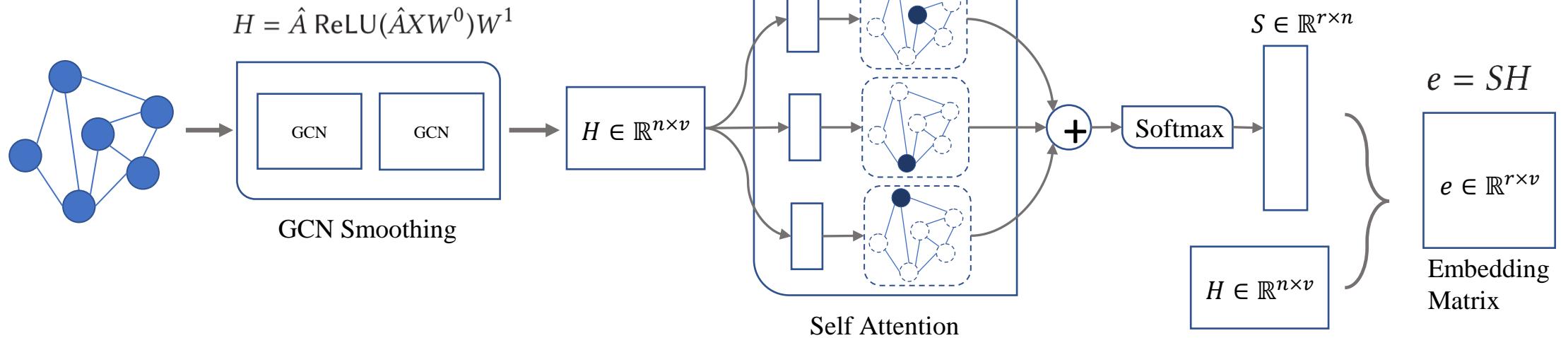
# Hierarchical Graph Classification

- **The Problem:** predicts the class label of graph instances in a hierarchical graph.
- **Challenges:**
  - How to represent the graphs with arbitrary size into a fixed-length vector?
  - How to incorporate the information of instance level and hierarchical level?



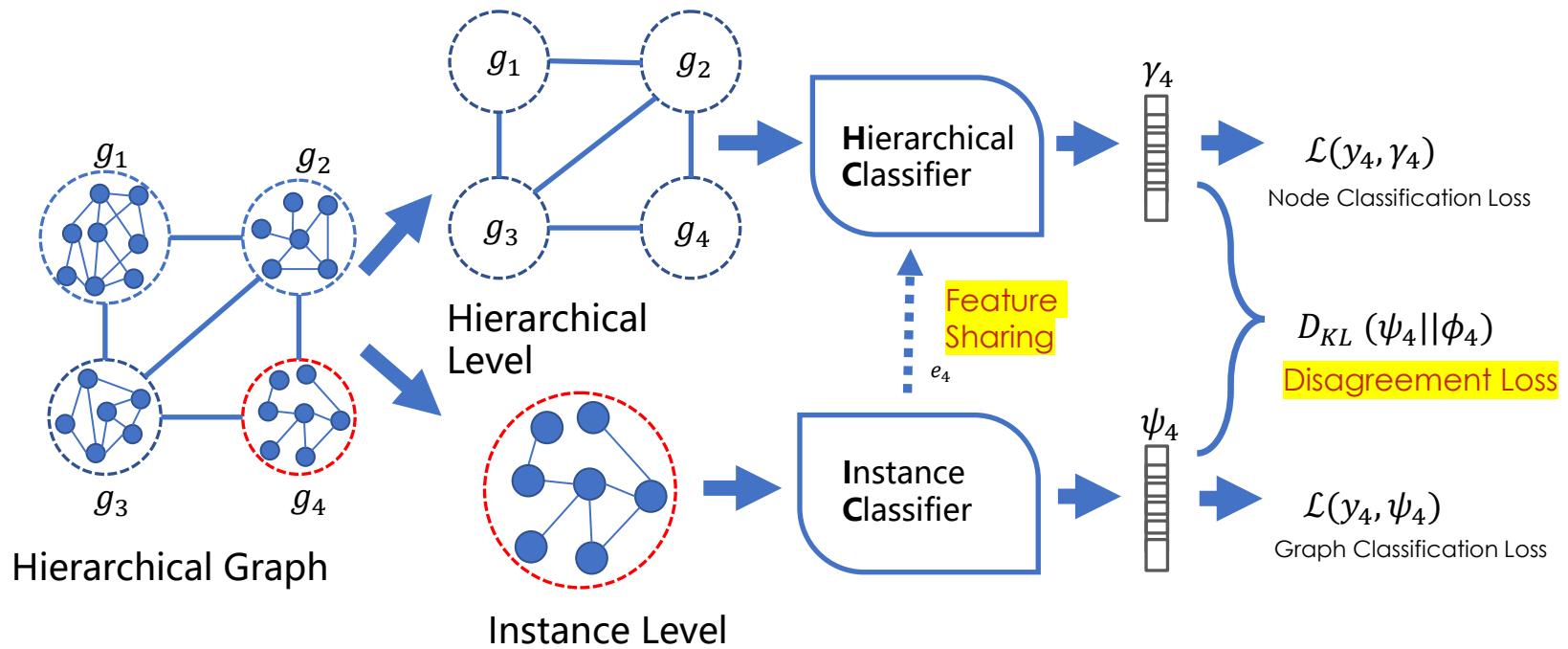
# Graph Instance Level: Self-Attentive Graph Embedding

- How to represent the graphs with arbitrary size into a fixed-length vector?
- Graph representation learning at different level:
  - Node Level:  $G(V, E) \rightarrow H^{n \times v}$
  - Graph Level:  $G(V, E) \rightarrow e^v$
- **SAGE: Self-Attentive Graph Embedding**
  - Size invariance ---- Self-attention
  - Permutation invariance ---- GCN Smoothing
  - Node importance ---- Self-attention
- Self-attention  $S$  :  $r$  opinions about node importance.



# The Unified Model

- How to incorporate the information of instance level and hierarchical level?
  - **Instance Level Model** : Graph Level Learning (SEGA)
  - **Hierarchical Level Model**: Node Level Learning (GCN)
- **Feature Sharing**: Concatenate the output of SEGA to the input of GCN.
- **Disagreement Loss**: The disagreement between instance classifier and hierarchical classifier should be minimized.



The overall loss:  $\min \zeta(G_l) + \xi(G_u)$ ,

The supervised loss:

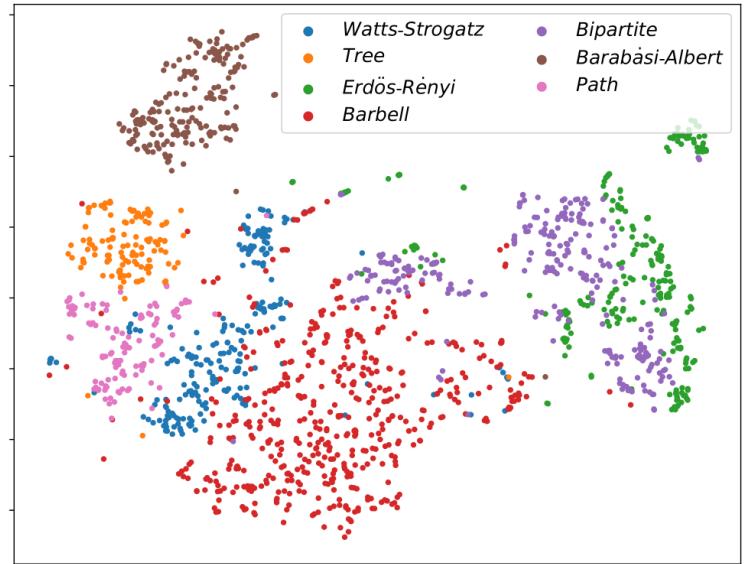
$$\zeta(G_l) = \sum_{g_i \in G_l} (\mathcal{L}(y_i, \psi_i) + \mathcal{L}(y_i, \gamma_i)),$$

The disagreement loss:

$$\xi(G_u) = \sum_{g_i \in G_u} D_{KL}(\gamma_i || \psi_i),$$

# Performance of SAGE

- Performance of synthesized graphs
- Performance on the protein classification task:



Two-dimensional visualization of graph embeddings generated from the synthesized graph instances using SAGE.

**Table 1: Statistics of PROTEINS and D&D**

	PROTEINS	D&D
Max number of nodes	620	5748
Avg number of nodes	39.06	284.32
Number of graphs	1113	1178

**Table 2: Accuracy of different classifiers**

Approach	PROTEINS	D&D
SP	$75.07 \pm 0.54\%$	-
RW	$74.22 \pm 0.42\%$	-
GK	$71.67 \pm 0.55\%$	$78.45 \pm 0.26\%$
WL	$72.92 \pm 0.56\%$	$77.95 \pm 0.70\%$
PSCN	$75.89 \pm 2.76\%$	$77.12 \pm 2.41\%$
graph2vec	$73.30 \pm 2.05\%$	-
SAGE	$77.26 \pm 2.28\%$	$80.88 \pm 2.33\%$

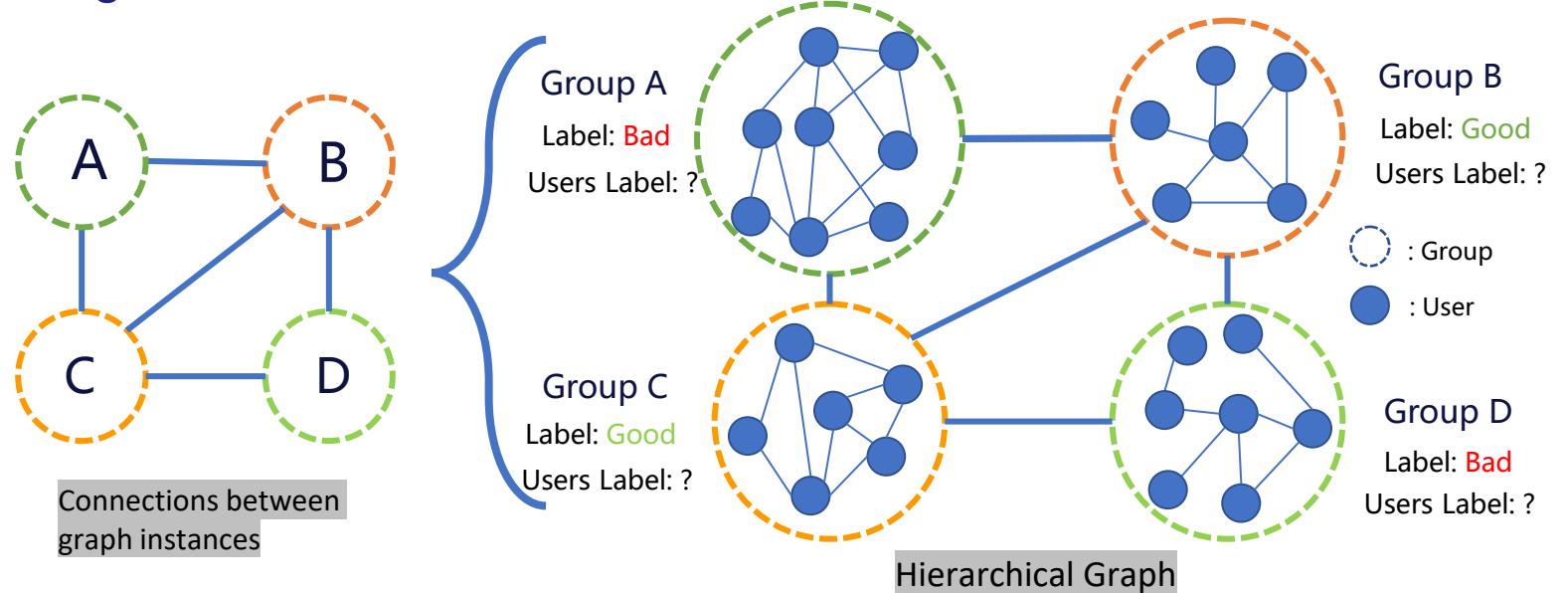
# GNN in Social Networks

- "Semi-supervised graph classification: A hierarchical graph perspective." **WWW 2019**
- "Inverse Graph Identification: Can We Identify Node Labels Given Graph Labels?" **arXiv 2020**

# Inverse Graph Identification (IGI)

## Inverse Graph Identification (IGI):

- **The problem:** infers the information or properties of nodes or even edges and sub-graphs in each graph instance.
- A node clustering case of IGI:

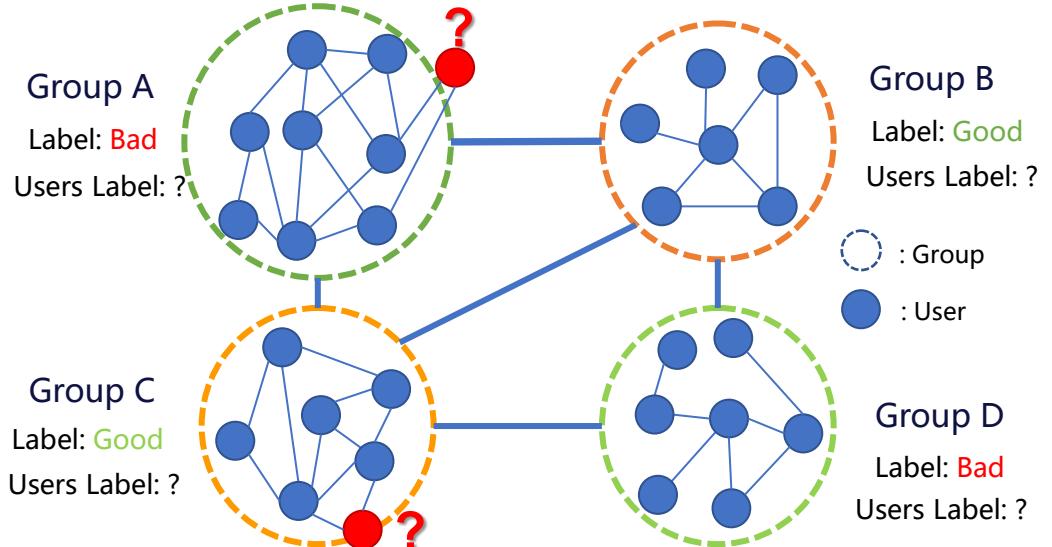
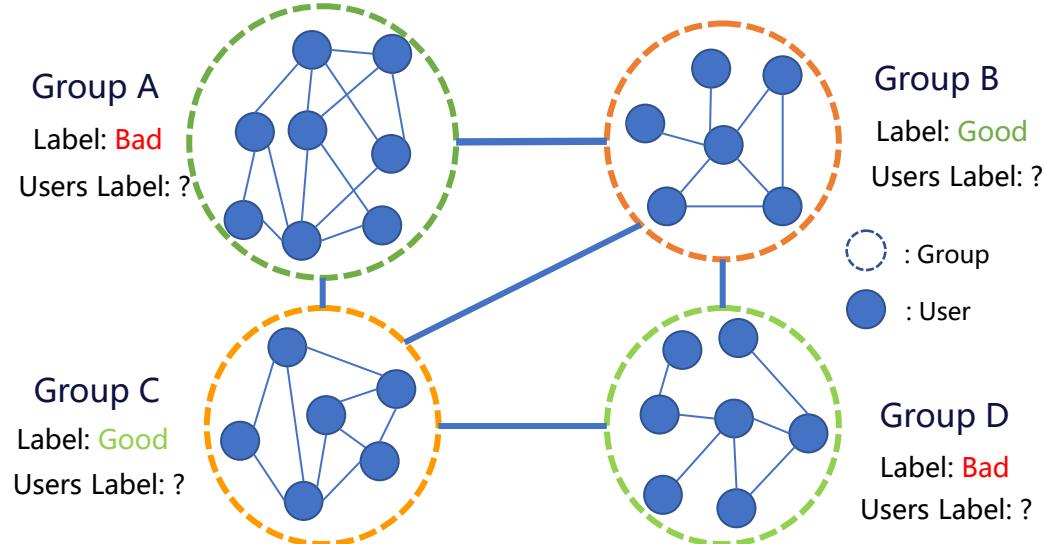


- E.g., social media, point cloud segmentation, etc.

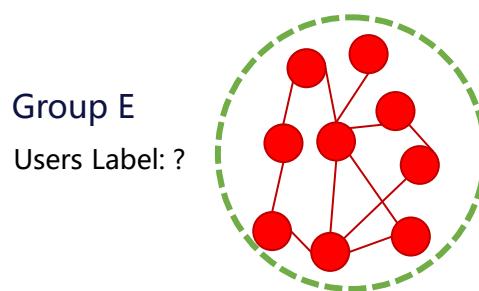


# IGI: Different Cases

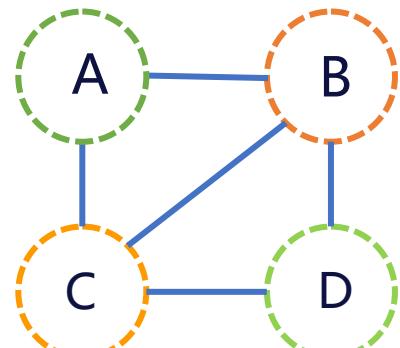
## Case 1: Infer existing node labels Case 2: Infer new node labels



## Case 3: Infer node labels of new graphs

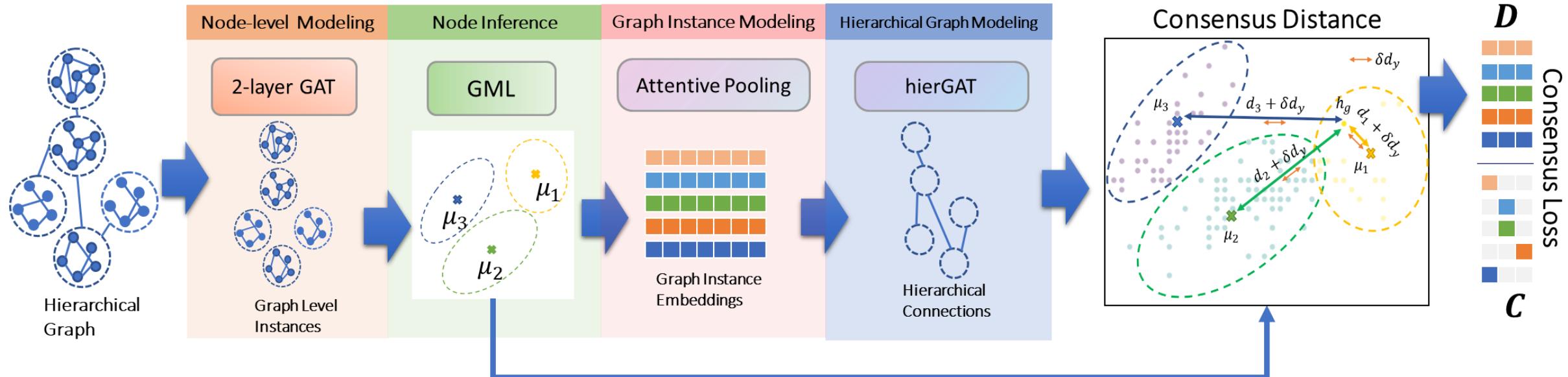


## Case 4: With hierarchical graph



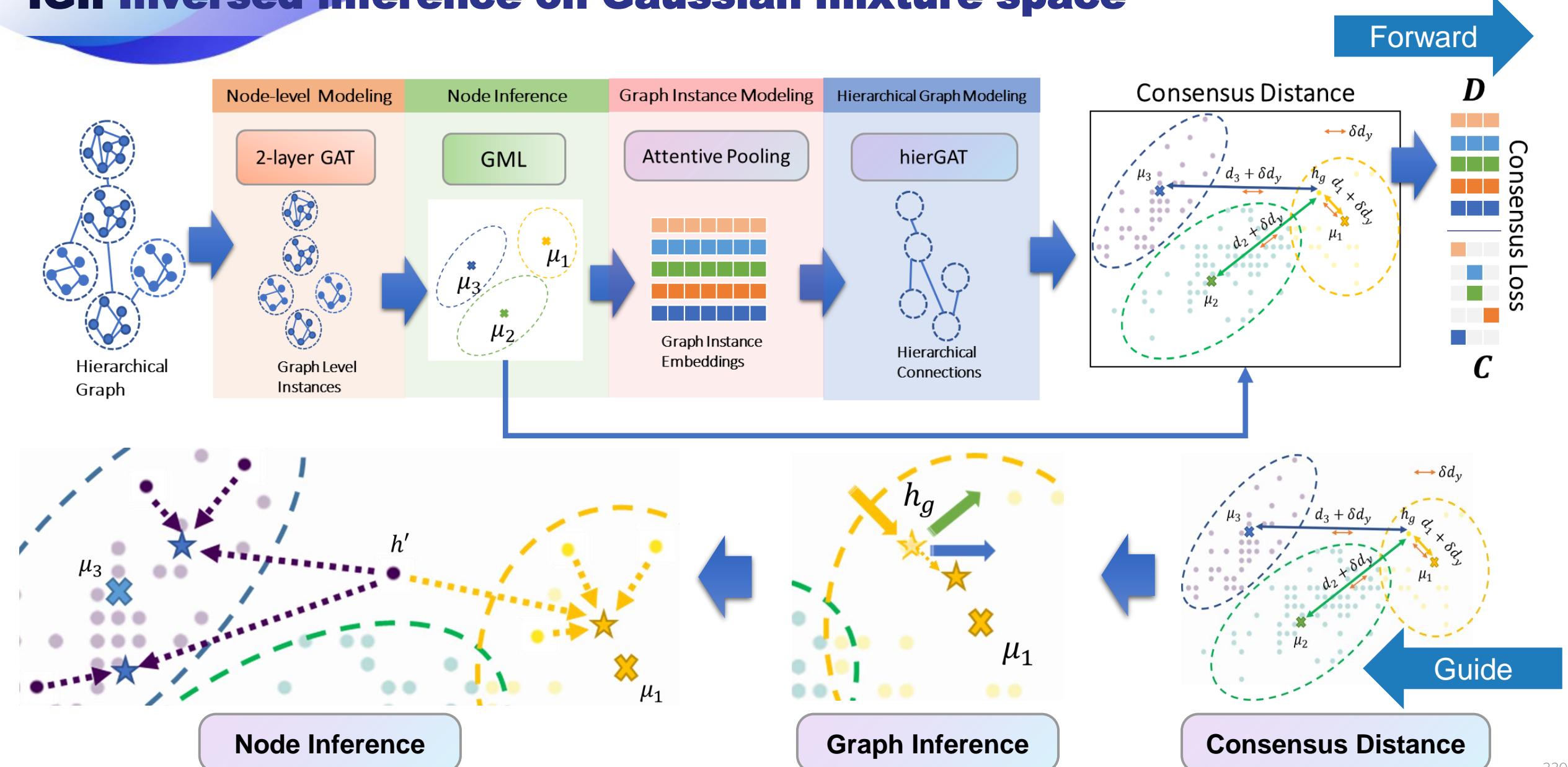
# IGI: Inverse Graph Identification on Hierarchical Graph

## Gaussian Mixture Graph Convolutional Networks (GMGCN):



- **Gaussian Mixture Layer (GML):**  $h' = \frac{1}{c} \sum_{c=1}^C \exp\left(-\frac{1}{2\sigma_c^2} (h - \mu_c)^2\right)$
- **Consensus Distance:**  $d_n = [\|h_g^n - \mu_1\|_2, \|h_g^n - \mu_2\|_2, \dots, \|h_g^n - \mu_c\|_2] \in \mathbb{R}^c$
- **Consensus Loss:**  $\mathcal{L}_{con} = \text{cross\_entropy}_n(\text{softmax}(-d_n - \delta \times \text{onehot}(c_n) \odot d_n), c_n)$
- **Node Inference:**  $z_i = \arg \max_c \text{similarity}(h'_i, \mu_c)$

# IGI: Inversed inference on Gaussian mixture space



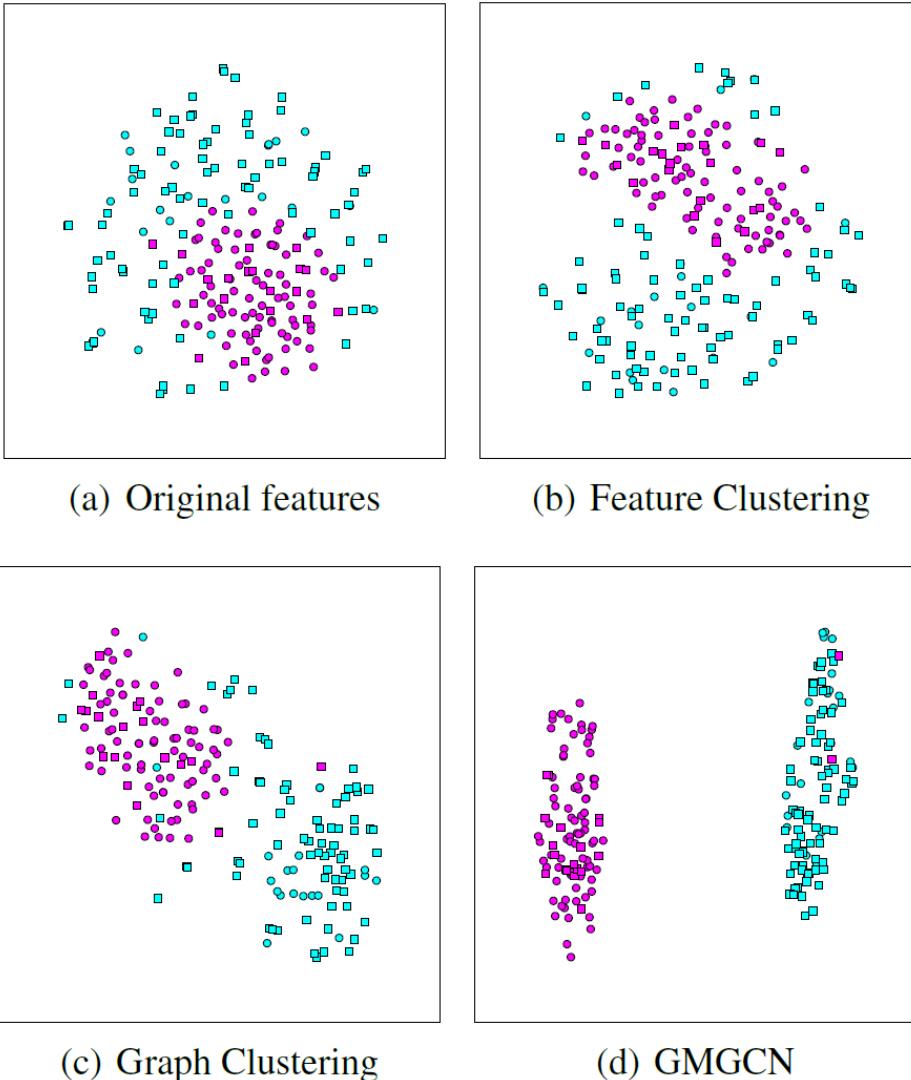


# IGI: Inverse Graph Identification on Hierarchical Graph

- Performance of synthesized graphs

methods	single graph					
	original nodes		new nodes		new graphs	
	NMI	ARI	NMI	ARI	NMI	ARI
Feature Clustering	49.46(1.89)	50.78(2.48)	50.11(2.07)	49.33(2.56)	47.48(2.48)	48.03(3.32)
Graph Clustering	80.65(8.18)	85.51(7.96)	85.85(8.47)	87.72(8.52)	78.72(10.47)	84.25(10.33)
MIL	76.47(0.00)	81.39(0.00)	75.42(0.00)	80.76(0.00)	75.77(0.00)	80.63(0.00)
ATTGCN	14.25(9.76)	12.45(11.64)	19.69(9.54)	13.44(11.38)	13.25(9.60)	12.22(11.10)
GMGCN	<b>92.26(0.22)</b>	<b>95.71(0.13)</b>	<b>94.70(0.33)</b>	<b>96.08(0.26)</b>	<b>92.98(0.24)</b>	<b>96.16(0.16)</b>

methods	hierarchical graph					
	original nodes		new nodes		new graphs	
	NMI	ARI	NMI	ARI	NMI	ARI
ATTGCN	19.09(12.34)	17.48(13.52)	24.64(12.34)	19.09(13.23)	20.54(11.92)	20.00(13.37)
GMGCN	<b>92.88(0.31)</b>	<b>96.05(0.17)</b>	<b>95.16(0.33)</b>	<b>96.42(0.25)</b>	<b>96.36(0.29)</b>	<b>96.39(0.16)</b>





# IGI: Inverse Graph Identification on Hierarchical Graph

- Performance of social media datasets (PHEME)

	Charlie Hebdo	Ferguson	Germanwings Crash	Ottawa Shooting	Sydney Siege
Method	NMI	NMI	NMI	NMI	NMI
Feature Clustering	25.46(0.00)	24.56(0.00)	48.84(0.00)	<b>34.56(0.00)</b>	19.45(0.00)
Graph Clustering	3.26(2.26)	1.08(0.71)	0.82(0.40)	4.49(3.07)	4.93(1.92)
MIL	5.69(0.00)	4.08(0.00)	0.61(0.00)	0.60(0.00)	19.68(0.00)
Feature Clustering+feature augmentation	25.42(0.00)	24.57(0.00)	48.84(0.00)	34.33(0.00)	19.21(0.00)
Graph Clustering+feature augmentation	6.08(3.06)	3.45(2.90)	39.14(11.78)	6.22(4.15)	6.83(4.28)
MIL+feature augmentation	22.32(0.00)	0.14(0.00)	0.28(0.00)	0.11(0.00)	19.98(0.00)
GMGCN	<b>47.51(3.27)</b>	<b>48.35(4.08)</b>	<b>48.85(2.14)</b>	32.58(3.63)	<b>41.00(3.93)</b>
Method	ARI	ARI	ARI	ARI	ARI
Feature Clustering	23.46(0.00)	23.23(0.00)	47.29(0.00)	32.08(0.00)	16.45(0.00)
Graph Clustering	6.56(6.30)	4.04(2.24)	3.71(1.93)	10.66(7.68)	11.78(5.57)
MIL	9.38(0.00)	13.93(0.00)	1.10(0.00)	2.09(0.00)	<b>39.51(0.00)</b>
Feature Clustering+feature augmentation	23.41(0.00)	23.23(0.00)	48.29(0.00)	31.85(0.00)	16.20(0.00)
Graph Clustering+feature augmentation	10.11(6.02)	9.15(6.51)	10.12(4.47)	8.91(7.99)	12.65(6.99)
MIL+feature augmentation	42.51(0.00)	1.08(0.00)	0.04(0.00)	0.04(0.00)	33.71(0.00)
GMGCN	<b>52.38(3.15)</b>	<b>55.26(3.21)</b>	<b>54.95(1.02)</b>	<b>37.51(2.23)</b>	37.79(2.18)



Tencent  
AI Lab



# Applications

## GNN in Medical Imaging

# GNN in Medical Imaging

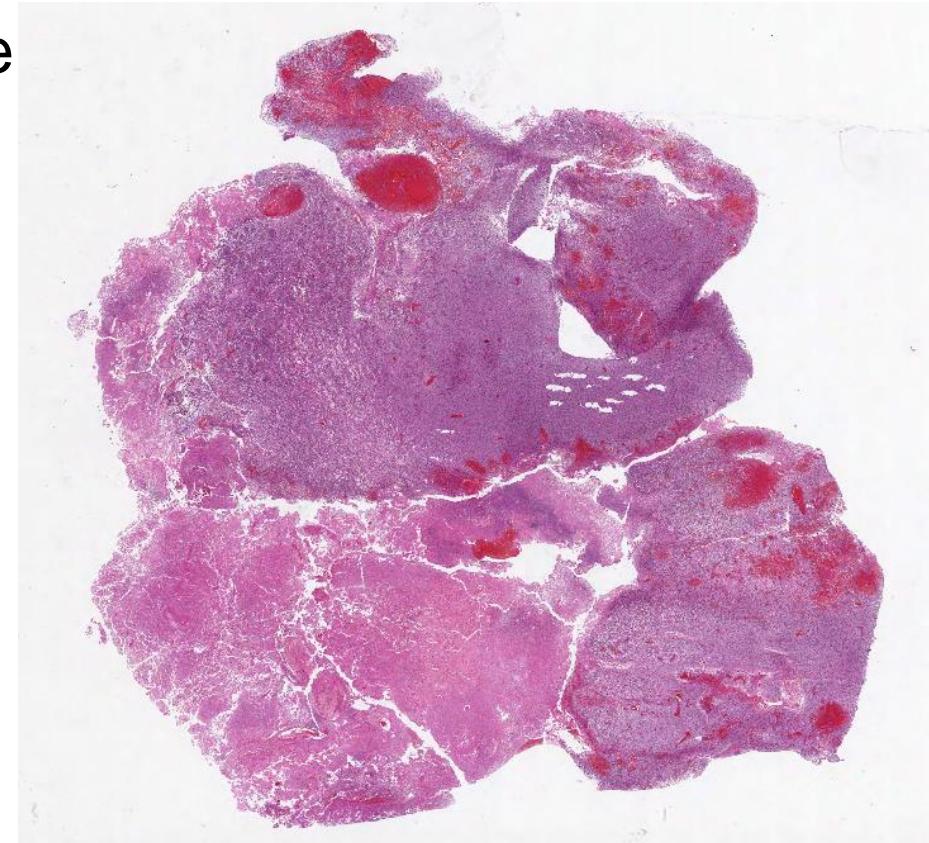
- "Graph CNN for Survival Analysis on Whole Slide Pathological Images", **MICCAI 2018**
- "Graph Convolutional Nets for Tool Presence Detection in Surgical Videos", **IPMI 2019**
- "Graph Attention Multi-instance Learning for Accurate Colorectal Cancer Staging", **MICCAI 2020**

## • Survival Prediction

- Predict the risk of a certain event occurs.
- Event: part failure, drug adverse reaction or death.
- Application: provides suggestion for clinical inte

## • Whole Slide Images

- Large: single WSI size  $>0.5$  GB.
- Complicated: millions of cells.
- Combine local and global features.



- **Cox**  $\lambda(t|X_i) = \lambda_0(t) \exp(\beta_1 X_{i1} + \cdots + \beta_p X_{ip}) = \lambda_0(t) \exp(X_i \cdot \beta)$

- **Partial likelihood for event happens on subject  $i$ :**

$$L_i(\beta) = \frac{\lambda(Y_i|X_i)}{\sum_{j:Y_j \geq Y_i} \lambda(Y_j|X_j)} = \frac{\cancel{\lambda_0(Y_i)}\theta_i}{\sum_{j:Y_j \geq Y_i} \cancel{\lambda_0(Y_j)}\theta_j} = \frac{\theta_i}{\sum_{j:Y_j \geq Y_i} \theta_j}$$

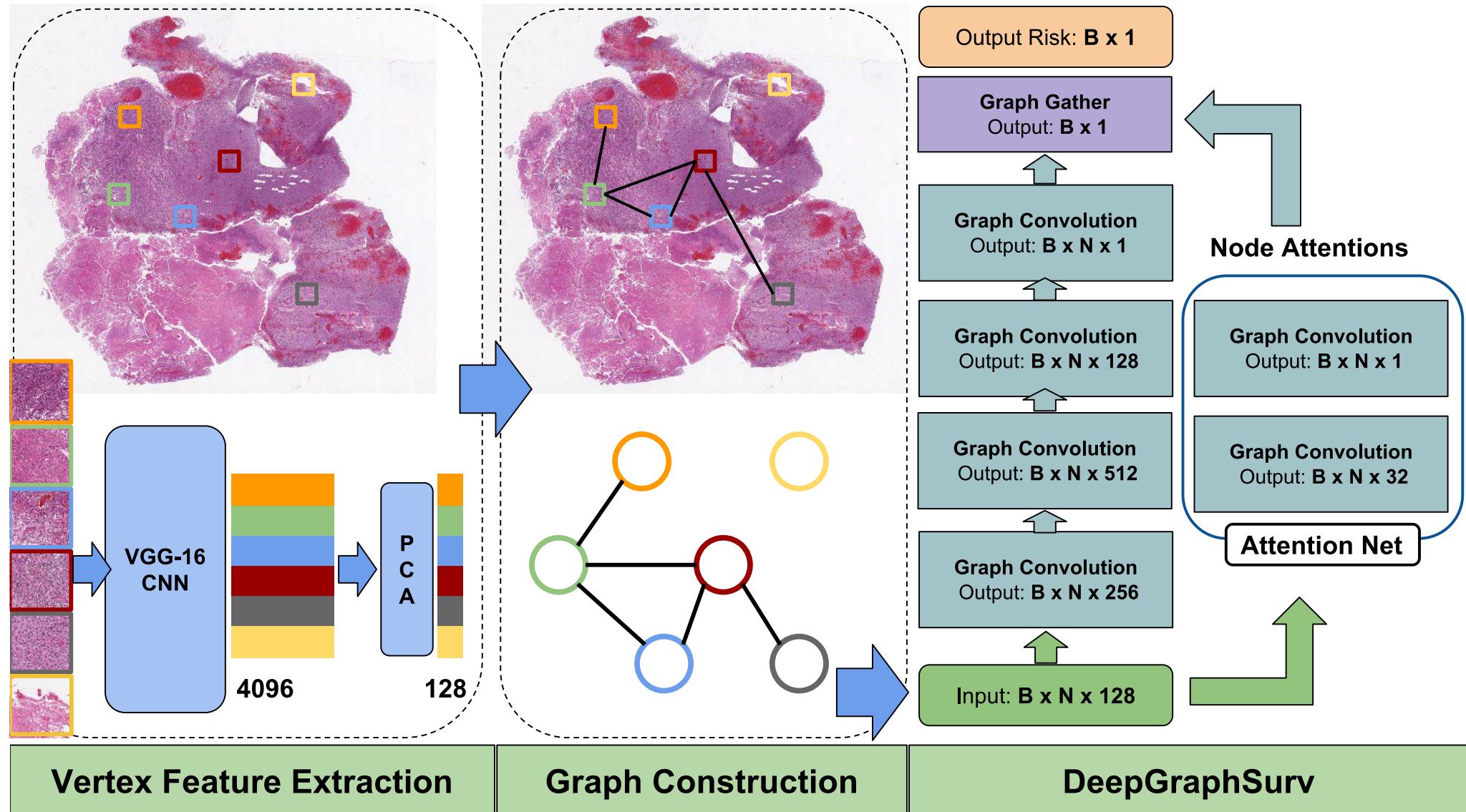
$$\theta_j = \exp(X_j \cdot \beta)$$

where,  $Y$  is the observation time.

- **Join likelihood of all subjects:**  $L(\beta) = \prod_{i:C_i=1} L_i(\beta)$

- **Log likelihood as object function:**

$$\ell(\beta) = \sum_{i:C_i=1} \left( X_i \cdot \beta - \log \sum_{j:Y_j \geq Y_i} \theta_j \right)$$



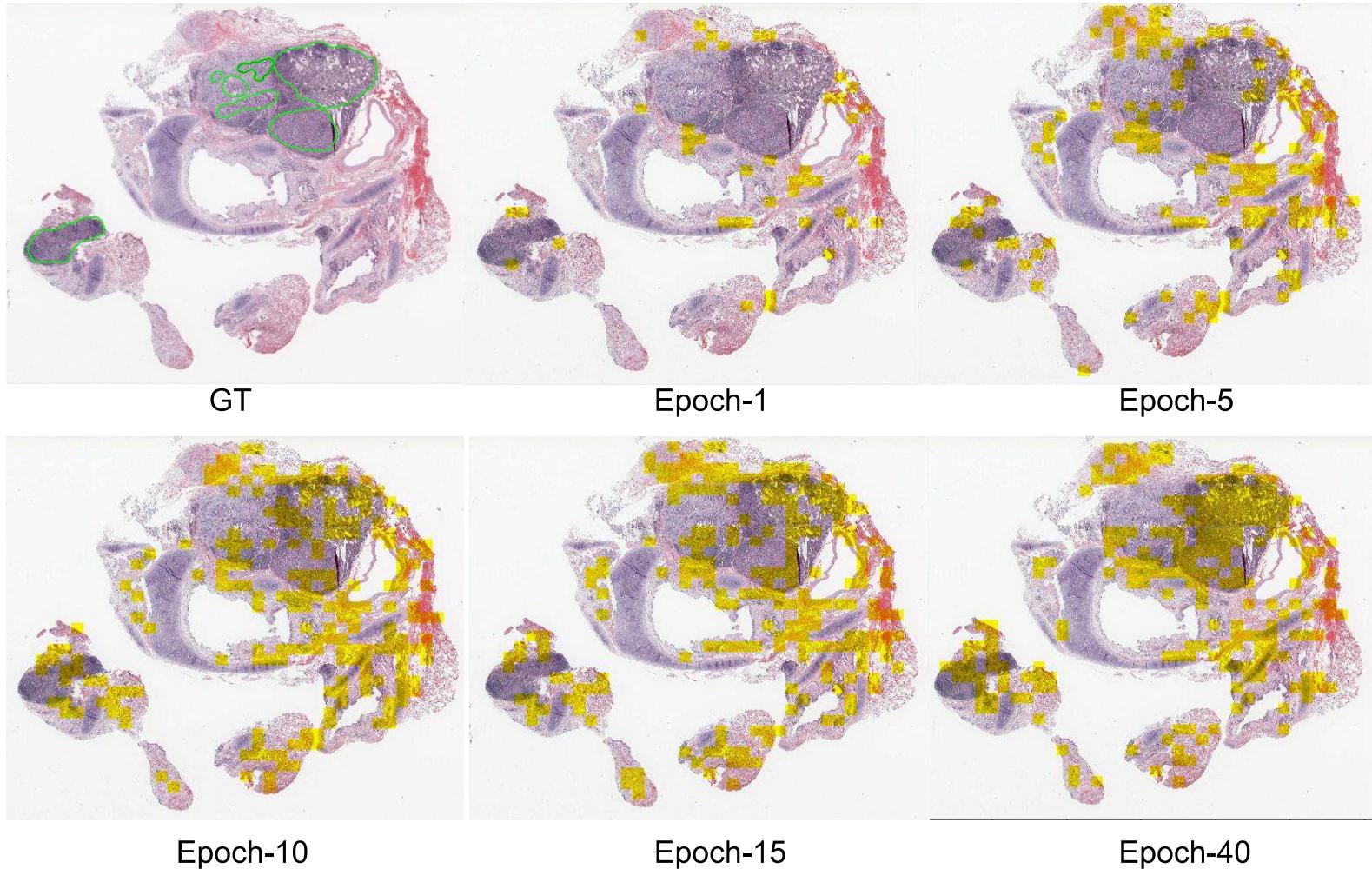
- **Pathological Images and Patient Survival Time and Label**

- TCGA, The Cancer Genome Atlas
- NLST, National Lung Screening Trials

Database	Cancer Subtype	No. Patient	No. WSI	Quality	Avg. Size
TCGA	LUSC	463	535	medium	0.72 GB
TCGA	GBM	365	491	low	0.50 GB
NLST	ADC & SCC	263	425	high	0.74 GB

- **Evaluation Metrics**- C-index: the fraction of all pairs of patients whose predicted survival times are correctly ordered.

- **Yellow regions: high attention values**
  - High attention patches : values > 0.9 (attention values (0, 1))

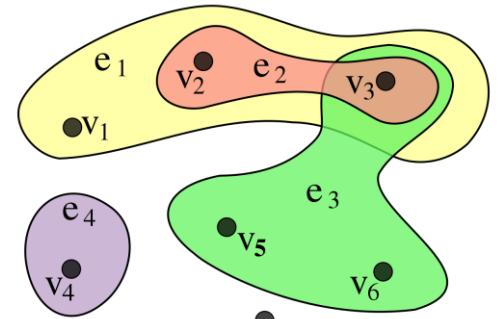
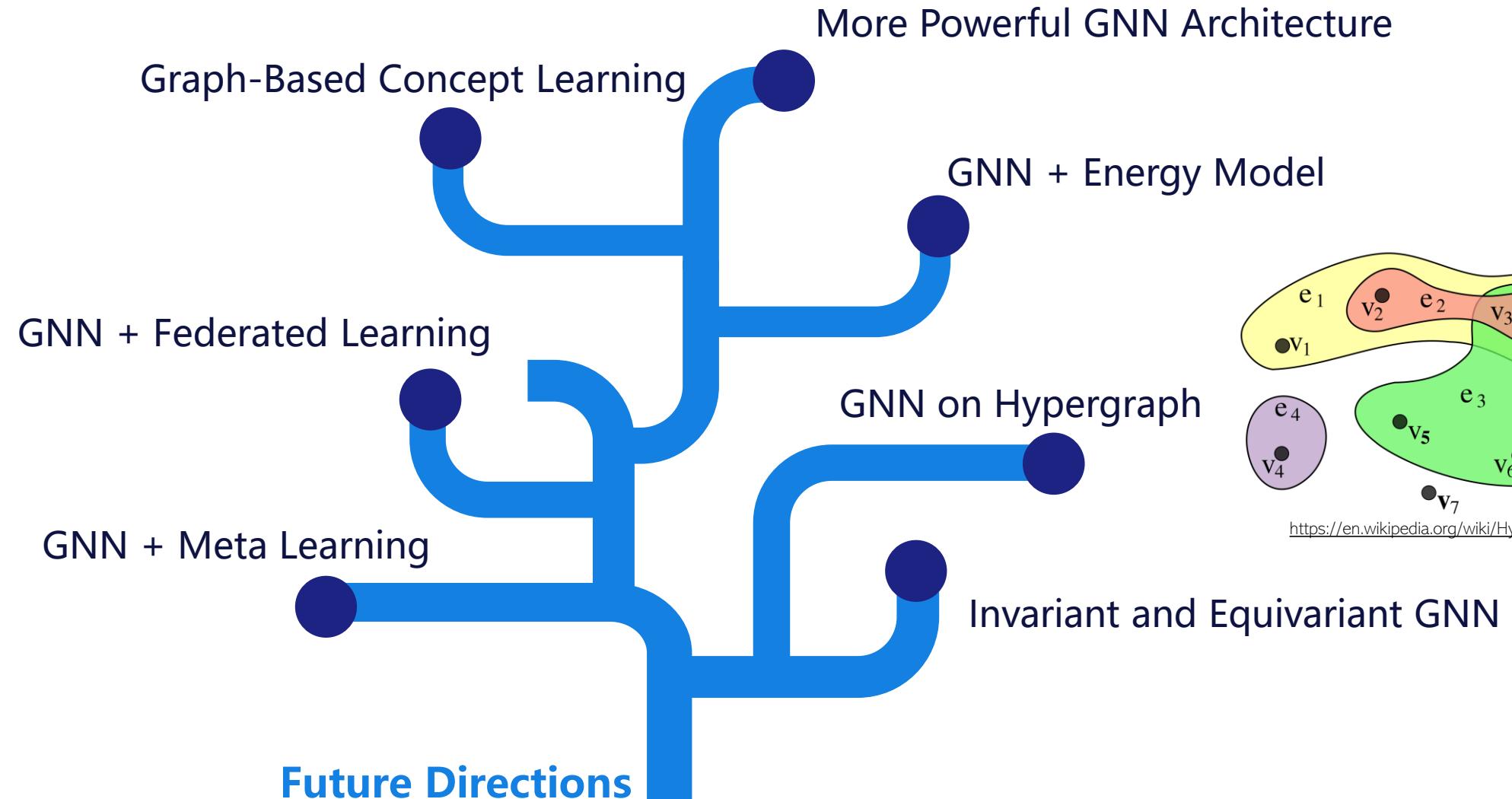




Model	LUSC	GBM	NLST
LASSO-Cox [19]	0.5280	0.5574	0.4738
LASSO-Cox*	<b>0.5663</b>	0.5165	<b>0.5663</b>
BoostCI [17]	0.5633	0.5543	0.5705
BoostCI*	<b>0.5800</b>	0.5130	<b>0.5716</b>
EnCox [20]	0.5216	0.5597	0.4883
EnCox*	<b>0.5740</b>	0.5231	<b>0.5742</b>
RSF [12]	0.5066	0.5570	0.5964
RSF*	<b>0.5492</b>	0.5193	0.5491
MTLSA [16]	0.5386	0.5787	0.6042
MTLSA*	0.5247	0.5630	0.5573
WSISA [21]	0.6380	0.5760	0.6539
GCN-Cox [8]	0.6280	0.5901	0.6845
<b>DeepGraphSurv</b>	<b>0.6606</b>	<b>0.6215</b>	<b>0.7066</b>

\* Use our graph features for the survival model.

# Future Directions



<https://en.wikipedia.org/wiki/Hypergraph>



Tencent  
AI Lab



# Bibliography



# Bibliography

- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. NeurIPS, 2018.
- Abu-El-Haija, Sami, et al. "MixHop: Higher-Order Graph Convolutional Architectures via Sparsified Neighborhood Mixing." International Conference on Machine Learning. 2019.
- Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. ICLR 2017.
- Atwood, James, and Don Towsley. "Diffusion-convolutional neural networks." Advances in neural information processing systems. 2016.
- Belghazi et al. Mine: mutual information neural estimation. ICML, 2018.
- Bian, Tian, et al. "Inverse Graph Identification: Can We Identify Node Labels Given Graph Labels?." arXiv preprint arXiv:2007.05970 (2020).
- Bian, Tian, et al. "Rumor Detection on Social Media with Bi-Directional Graph Convolutional Networks." AAAI 2020.
- Bojchevski, Aleksandar, and Stephan Günnemann. "Adversarial attacks on node embeddings via graph poisoning." International Conference on Machine Learning. PMLR, 2019.
- Bruna, Joan, et al. "Spectral networks and locally connected networks on graphs." arXiv preprint arXiv:1312.6203 (2013).
- Chami, Ines, et al. "Hyperbolic graph convolutional neural networks." Advances in neural information processing systems. 2019.
- Chang, Heng, et al. "A restricted black-box adversarial framework towards attacking graph embedding models." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 34. No. 04. 2020.
- Chang, Heng, et al. "Spectral Graph Attention Network." arXiv preprint arXiv:2003.07450 (2020).
- Chen et al. On the equivalence between graph isomorphism testing and function approximation with GNNs. NeurIPS, 2019.
- Chen, Jianfei, Jun Zhu, and Le Song. "Stochastic Training of Graph Convolutional Networks with Variance Reduction." International Conference on Machine Learning. 2018.
- Chen, Jie, Tengfei Ma, and Cao Xiao. "FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling." International Conference on Learning Representations. 2018.
- Chen, Yuzhao, et al. "On Self-Distilling Graph Neural Network." arXiv preprint arXiv:2011.02255 (2020).
- Chiang, Wei-Lin, et al. "Cluster-GCN: An efficient algorithm for training deep and large graph convolutional networks." Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems, 2(4):303–314, 1989.
- Dai, Hanjun, et al. "Adversarial attack on graph structured data." International conference on machine learning. PMLR, 2018.
- Defferrard, Michaël, Xavier Bresson, and Pierre Vandergheynst. "Convolutional neural networks on graphs with fast localized spectral filtering." Advances in neural information processing systems. 2016.
- Dehmamy et al. Understanding the Representation Power of Graph Neural Networks in Learning Graph Topology. NeurIPS, 2019.
- Durán & Niepert. Learning Graph Representations with Embedding Propagation. NeurIPS, 2017.
- Entezari, Negin, et al. "All you need is low (rank) defending against adversarial attacks on graphs." Proceedings of the 13th International Conference on Web Search and Data Mining. 2020.
- Fan, Wenqi, et al. "Graph neural networks for social recommendation." The World Wide Web Conference. 2019.
- Garg et al. Generalization and Representational Limits of Graph Neural Networks. ICML, 2020.

# Bibliography

- Gilmer, Justin, et al. "Neural Message Passing for Quantum Chemistry." ICML. 2017.
- Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." ICLR 2015.
- Gori, Marco, Gabriele Monfardini, and Franco Scarselli. "A new model for learning in graph domains." Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.. Vol. 2. IEEE, 2005.
- Gou, Jianping, et al. "Knowledge distillation: A survey." International Journal of Computer Vision (2021): 1-31.
- Hamilton, Will, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large graphs." Advances in neural information processing systems. 2017.
- Hinton, G. E., & Salakhutdinov, R. R. Reducing the Dimensionality of Data with Neural Networks. Science, 2006.
- Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).
- Hjelm et al. Learning deep representations by mutual information estimation and maximization. ICLR, 2019.
- Hornik, K. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- Hu et al. Strategies for Pre-training Graph Neural Networks. ICLR, 2020.
- Huang, Wenbing, et al. "Adaptive sampling towards fast graph representation learning." Advances in neural information processing systems. 2018.
- Jin, Ming, et al. "Power up! Robust Graph Convolutional Network via Graph Powering." AAAI 2021
- Jin, Wei, et al. "Graph structure learning for robust graph neural networks." Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2020.
- Jin, Wengong, Regina Barzilay, and Tommi Jaakkola. "Multi-objective molecule generation using interpretable substructures." International Conference on Machine Learning. PMLR, 2020.
- Junchi, Yu, et al. "Graph Information Bottleneck for Subgraph Recognition.", ICLR 2021.
- Kipf & Welling. Variational Graph Auto-Encoders. arXiv, 2016.
- Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907 (2016).
- Klicpera, Johannes, Aleksandar Bojchevski, and Stephan Günnemann. "Predict then Propagate: Graph Neural Networks meet Personalized PageRank." International Conference on Learning Representations. 2018.
- Li, Jia, et al. "Adversarial attack on community detection by hiding individuals." Proceedings of The Web Conference 2020. 2020.
- Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective." The World Wide Web Conference. 2019.
- Li, Yujia, et al. "Gated graph sequence neural networks." arXiv preprint arXiv:1511.05493 (2015).
- Liao, Renjie, et al. "LanczosNet: Multi-Scale Deep Graph Convolutional Networks." International Conference on Learning Representations. 2018.
- Liu et al. N-Gram Graph: Simple Unsupervised Representation for Graphs, with Applications to Molecules. NeurIPS, 2019.
- Loukas, What graph neural networks cannot learn: depth vs width. ICLR, 2020.
- Ma, Yao, et al. "Graph convolutional networks with eigenpooling." Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.

# Bibliography

- Morris et al. Weisfeiler and Leman Go Neural Higher-order Graph Neural Networks. AAAI, 2019.
- Niepert, Mathias, Mohamed Ahmed, and Konstantin Kutzkov. "Learning convolutional neural networks for graphs." International conference on machine learning. 2016.
- Nowozin et al. f-gan: Training generative neural samplers using variational divergence minimization. NeurIPS, 2016.
- Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization. WWW, 2020.
- Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. KDD, 2020.
- Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data. arXiv, 2020.
- Sato et al. Approximation Ratios of Graph Neural Networks for Combinatorial Problems. NeurIPS, 2019.
- Scarselli, Franco, et al. "The graph neural network model." IEEE Transactions on Neural Networks 20.1 (2008): 61-80.
- Shervashidze et al. Weisfeiler-Lehman Graph Kernels. JMLR, 2011.
- Sperduti, Alessandro, and Antonina Starita. "Supervised neural networks for the classification of structures." IEEE Transactions on Neural Networks 8.3 (1997): 714-735.
- Sun et al. InfoGraph: Unsupervised and Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization. ICLR, 2020.
- Tai, Kai Sheng, Richard Socher, and Christopher D. Manning. "Improved semantic representations from tree-structured long short-term memory networks." arXiv preprint arXiv:1503.00075 (2015).
- Veličković et al. Deep Graph Infomax. ICLR, 2019.
- Veličković, Petar, et al. "Graph Attention Networks." International Conference on Learning Representations. 2018.
- Vincent et al. Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion. JMLR, 2010.
- Wang, X., & Gupta, A. Videos as Space-Time Region Graphs. ECCV, 2018.
- Xu et al. How Powerful Are Graph Neural Networks. ICLR, 2019.
- Xu, Bingbing, et al. "Graph Wavelet Neural Network." International Conference on Learning Representations. 2018.
- Xu, Kaidi, et al. "Topology attack and defense for graph neural networks: An optimization perspective." IJCAI 2019.
- Yan et al. Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. AAAI, 2018.
- Yang, Yiding, et al. "Distilling knowledge from graph convolutional networks." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- Ying, Rex, et al. "Graph convolutional neural networks for web-scale recommender systems." Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2018.
- Ying, Zhitao, et al. "Hierarchical graph representation learning with differentiable pooling." Advances in neural information processing systems. 2018.
- Yuan, Li, et al. "Revisiting knowledge distillation via label smoothing regularization." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.



# Bibliography

- Zeng, et al. Graph convolutional networks for temporal action localization. ICCV, 2019
- Zeng, Hanqing, et al. "GraphSAINT: Graph Sampling Based Inductive Learning Method." International Conference on Learning Representations. 2020.
- Zhang, Jiani, et al. "Gaan: Gated attention networks for learning on large and spatiotemporal graphs." arXiv preprint arXiv:1803.07294 (2018).
- Zhang, Xiang, and Marinka Zitnik. "Gnnguard: Defending graph neural networks against adversarial attacks." arXiv preprint arXiv:2006.08149 (2020).
- Zhu, Dingyuan, et al. "Robust graph convolutional networks against adversarial attacks." Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.
- Zügner, Daniel, Amir Akbarnejad, and Stephan Günnemann. "Adversarial attacks on neural networks for graph data." Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2018.
- Zügner, Daniel, and Stephan Günnemann. "Adversarial Attacks on Graph Neural Networks via Meta Learning." International Conference on Learning Representations. 2018.