

---

# Advanced Topics in Scalable Learning

CSE 6392 Lecture 2 Graph Neural Network

Junzhou Huang, Ph.D.

Department of Computer Science and Engineering

# Administration

---

- **Course CSE6392**

- What: Advanced Topics in Scalable Learning
- When: Friday 1:00 ~ 3:50pm
- Where: NH 111
- Who: Junzhou Huang (Office ERB 650) [jzhuang@uta.edu](mailto:jzhuang@uta.edu)
- Office Hour: FRIDAY 3:50 ~ 6:00pm and/or appointments
- Webpage: <http://ranger.uta.edu/~huang/teaching/CSE6392.htm>  
(Please check this page regularly)

- **Lecturer**

- PhD in CS from Rutgers, the State University of New Jersey
- Research areas: machine learning , computer vision, medical image analysis and bioinformatics

- **GTA**

- Saiyang Na (Office ERB 105B), [svn3892@mavs.uta.edu](mailto:svn3892@mavs.uta.edu)
- Office hours: Friday 10:00am ~ 12:00pm and/or appointments

# Assignment

---

- **Paper Selection**

- Each group has two members at most.
- Each group will select at least one paper from the following paper list and then be scheduled to present their selected papers in our class.
- You can choose any papers from the paper lists in the class
- Please talk to the lecturer if you prefer to select a paper out of the list
- The selected paper has to be confirmed by the second week
- GTA will set up the paper selection sheet
- Different groups will present different papers

# Grading

---

- **Distribution**

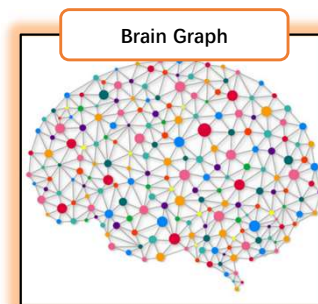
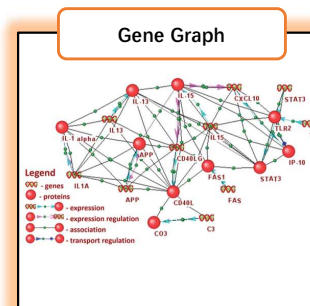
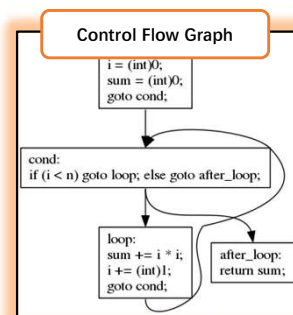
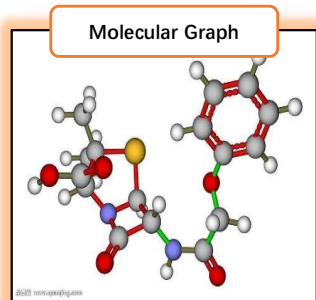
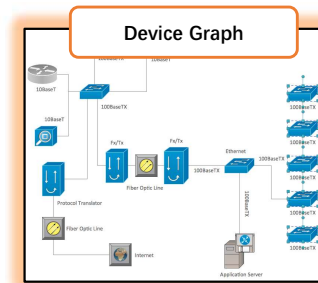
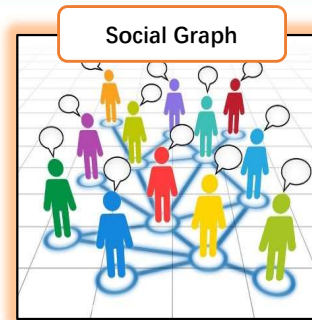
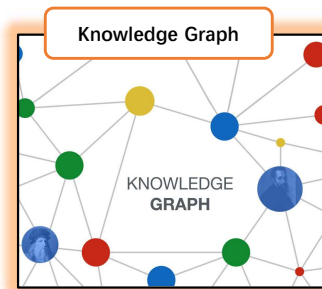
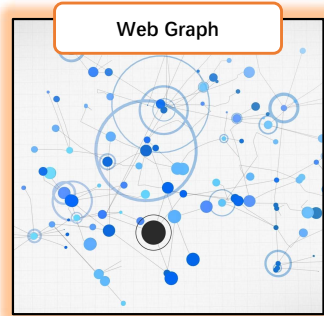
- 30% Paper Presentation
  - 30% Slide Preparation
  - 30% Questions & Answering
  - 10% Class Participation
- 
- 100%

- **Attention**

- No midterm or final exam for this course.
- Please read the selected paper and prepare the final presentation as early as possible
- This is research seminar course. Asking questions and discussion are highly encouraged



# Graph is Everywhere



# Graph is Important

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

### Social Recommendation

Fan W, Ma Y, Li Q, et al. 2019

### Survival Analysis

Li R, Yao J, Zhu X, et al. 2018

### Temporal Action Localization

Zeng et al. 2019

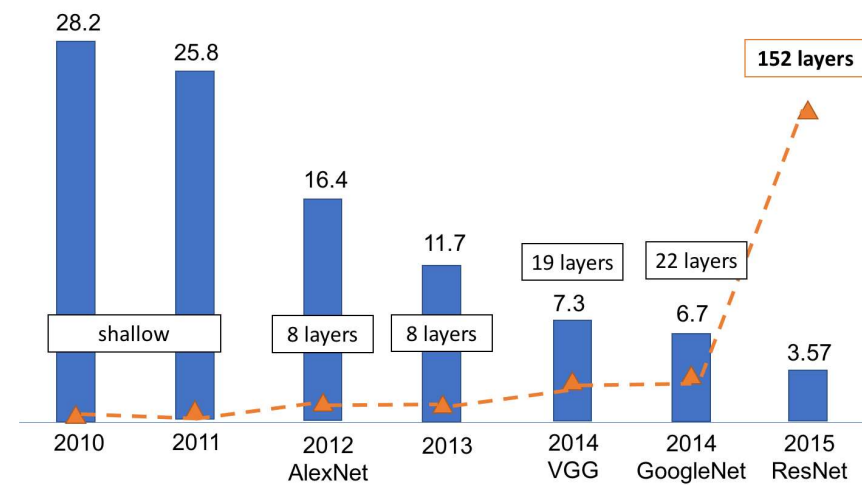
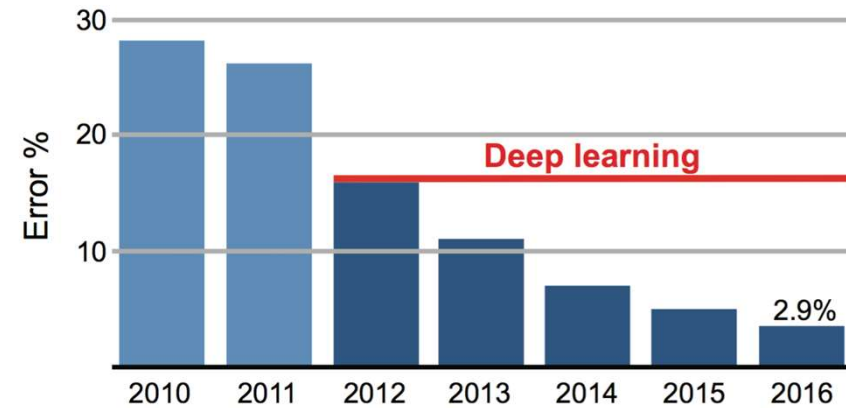
### Molecular Property Prediction

Property	Value
U0	-188.9736 (Hartree)
G	-188.9542 (Hartree)
H	-187.2784 (Hartree)

Lu C, Liu Q, Wang C, et al.

# The Power of Deep Learning

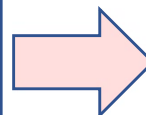
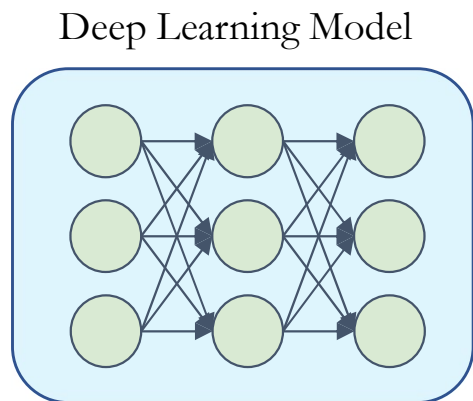
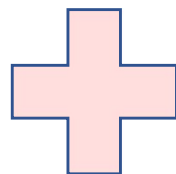
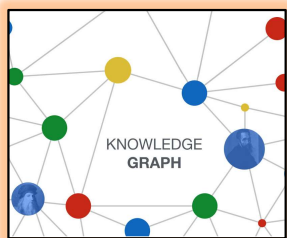
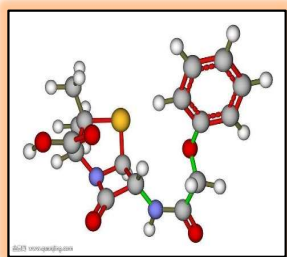
IMAGENET



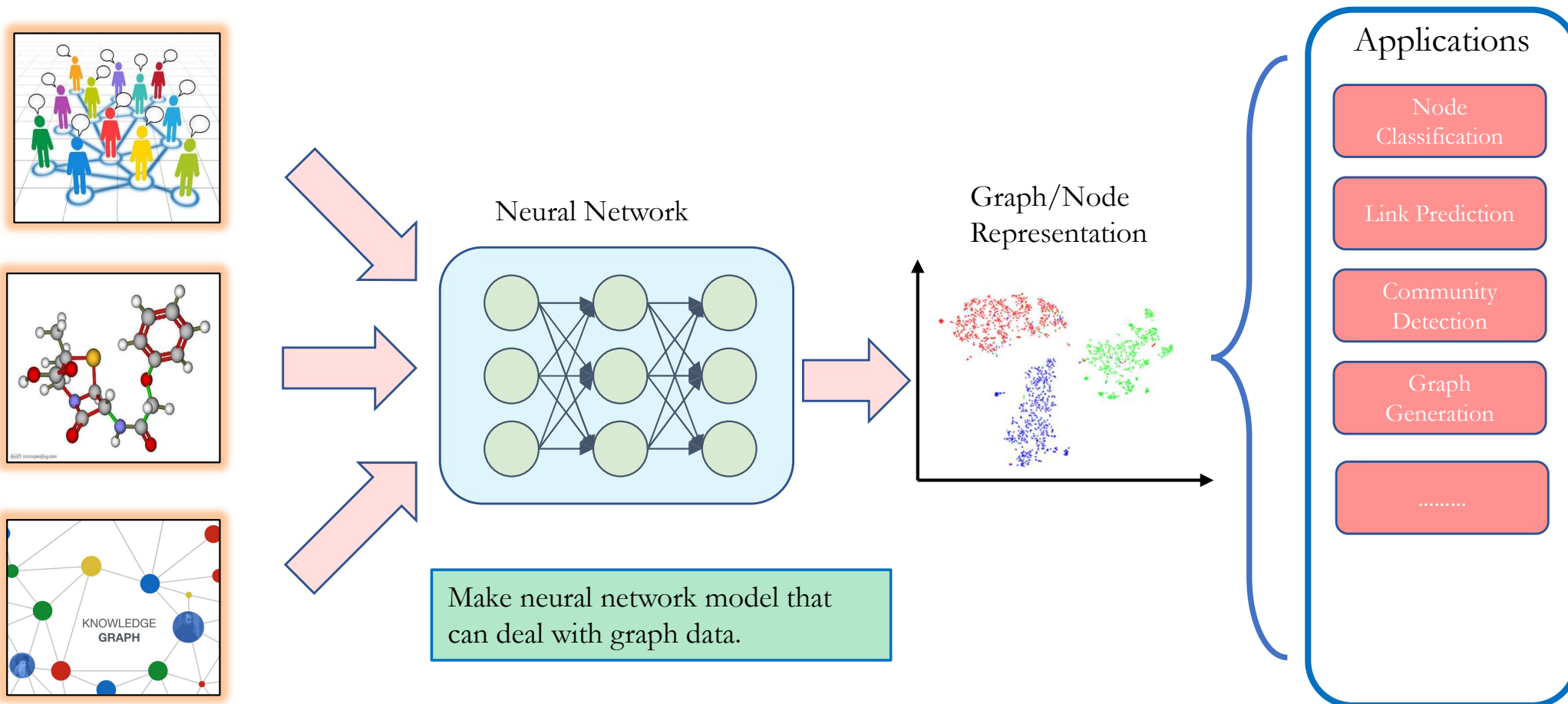


# Deep Learning + Graphs = ?

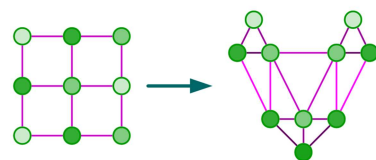
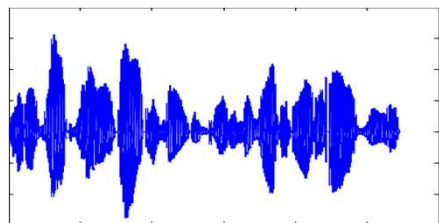
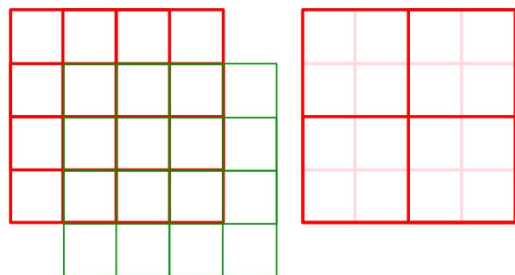
---



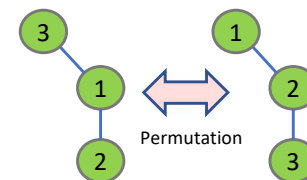
# Deep Graph Learning



# The Big Challenges

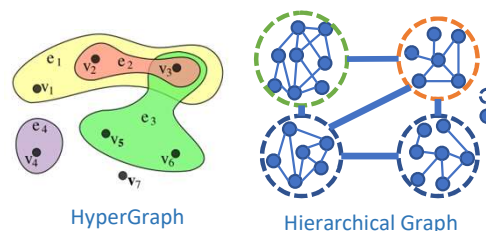


Irregular

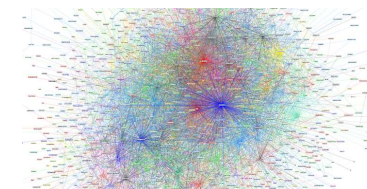


Permutation Invariant

VS



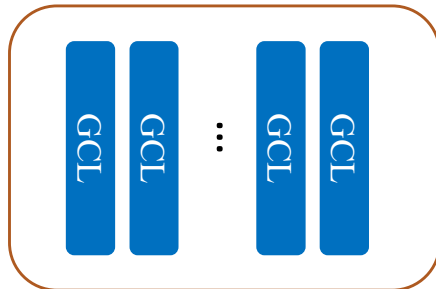
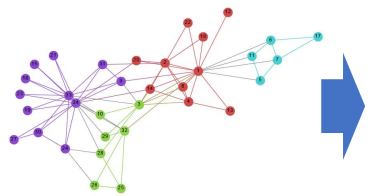
Complex variants



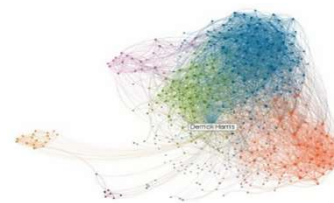
Large-scale instance

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

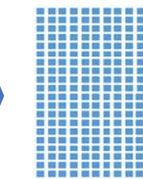
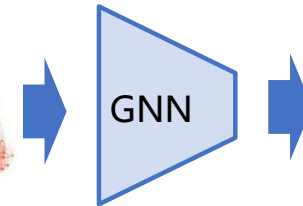
# The Research Questions



How to train the deep graph neural networks?

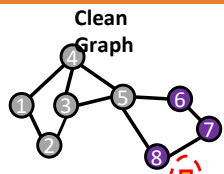


Big Graph

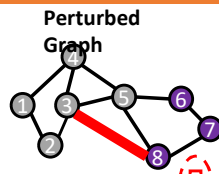


Node Representation

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

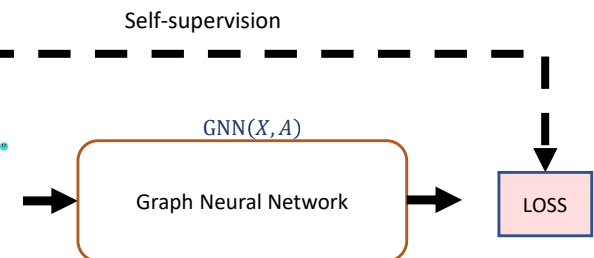
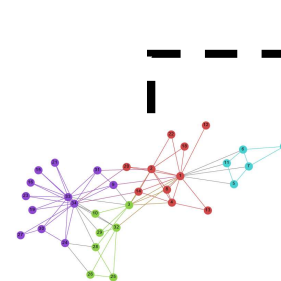


8 Predicted as:



8 Predicted as:

How to enhance the robustness of graph neural network?



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

# 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

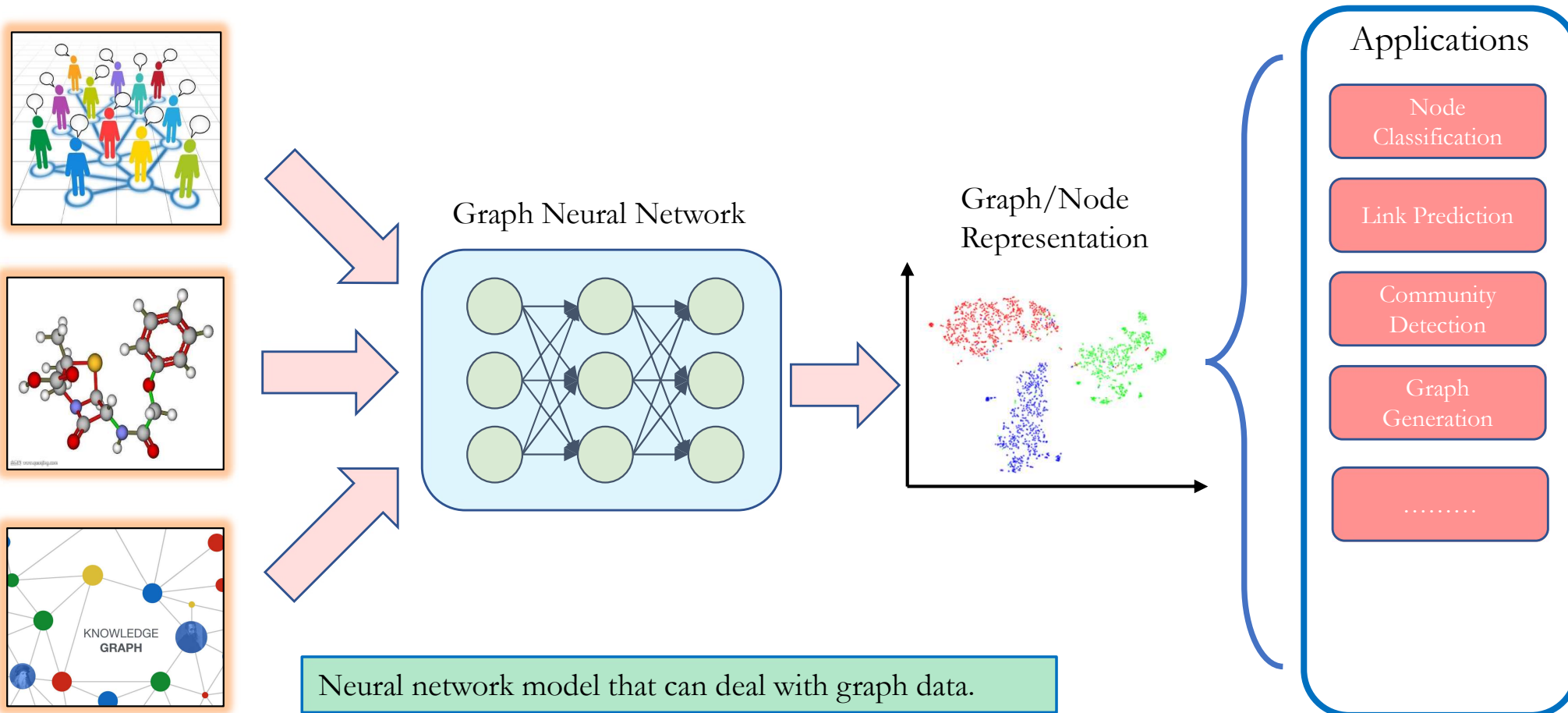
## Applications

- Social Networks
- Medical Imaging

---

# Preliminaries and Brief History of Graph Neural Networks

# What is the Graph Neural Network?



# Graph Neural Network is not a New Thing

---

Sperduti, Alessandro and Starita, Antonina. 1997

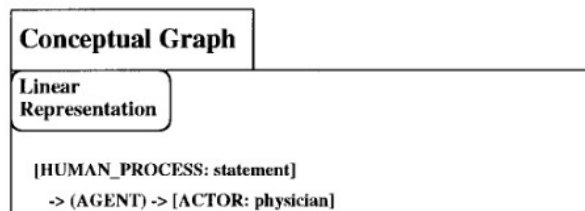
714

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.

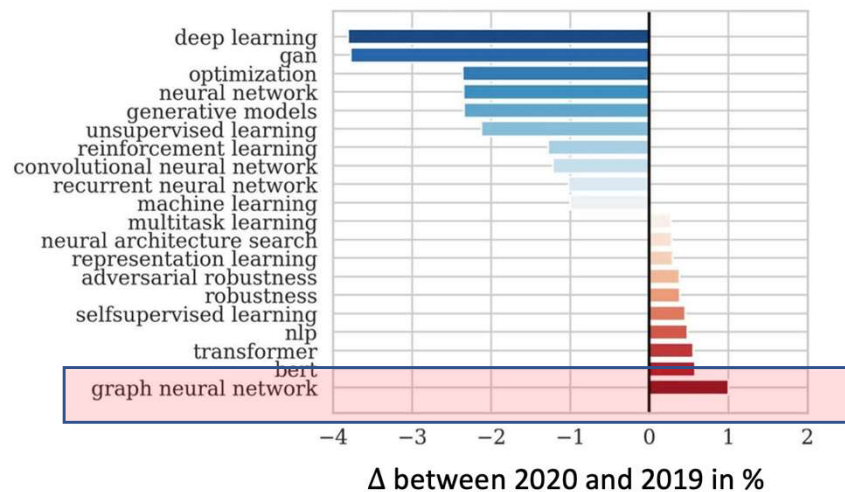


Sperduti, Alessandro, and Antonina Starita. "Supervised neural networks for the classification of structures."

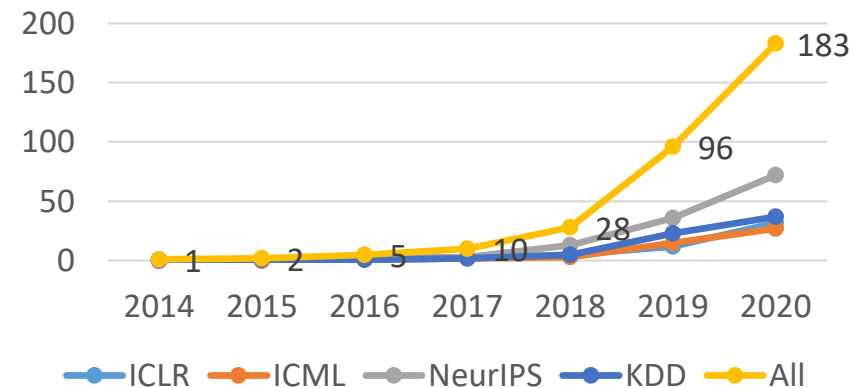


# A Rapidly Growing Area

ICLR 2020 submissions keyword statistics



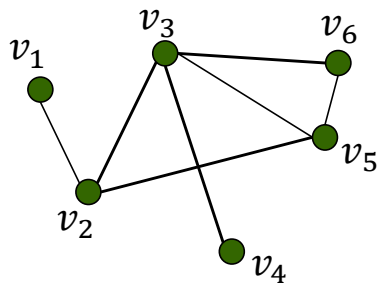
Number of GNN Papers



<https://github.com/shaohua0116/ICLR2020-OpenReviewData>

# 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:

$$\mathbf{X} \in \mathbb{R}^{n \times d}$$

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

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

**Laplacian matrix  $\mathbf{L} = \mathbf{D} - \mathbf{A}$**

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

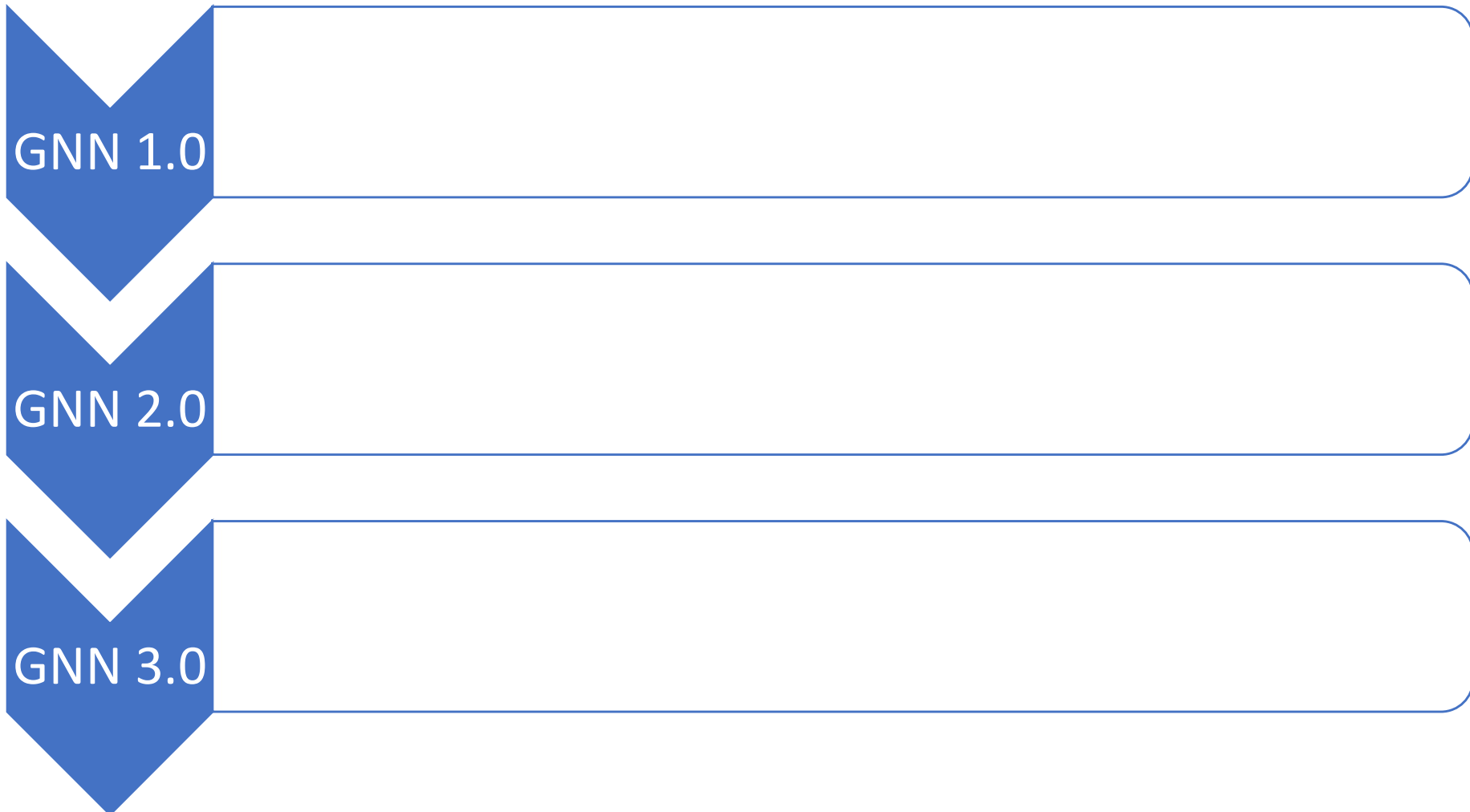
$$- \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{pmatrix} =$$

Laplacian matrix

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

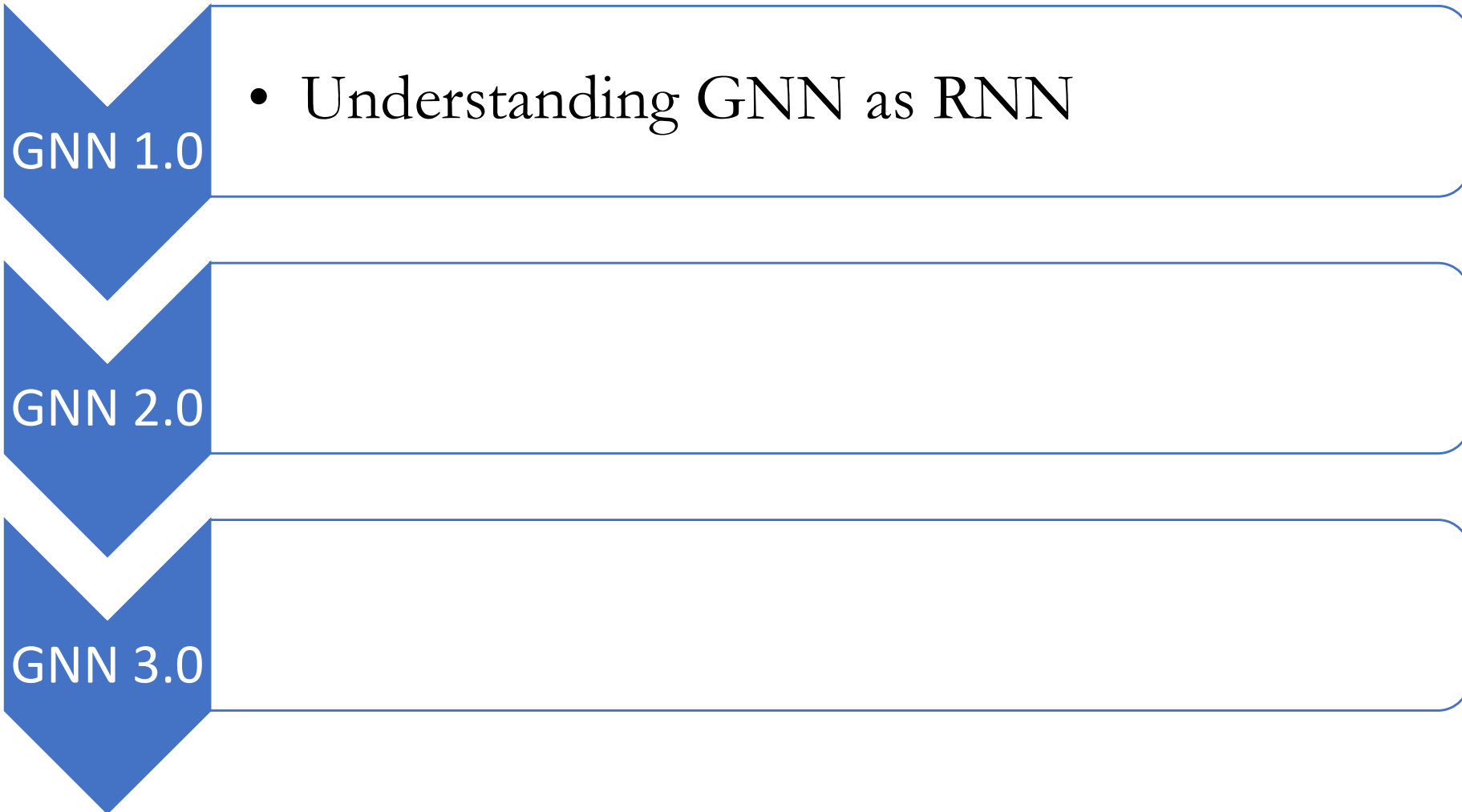
# The Model of Graph Neural Networks

---

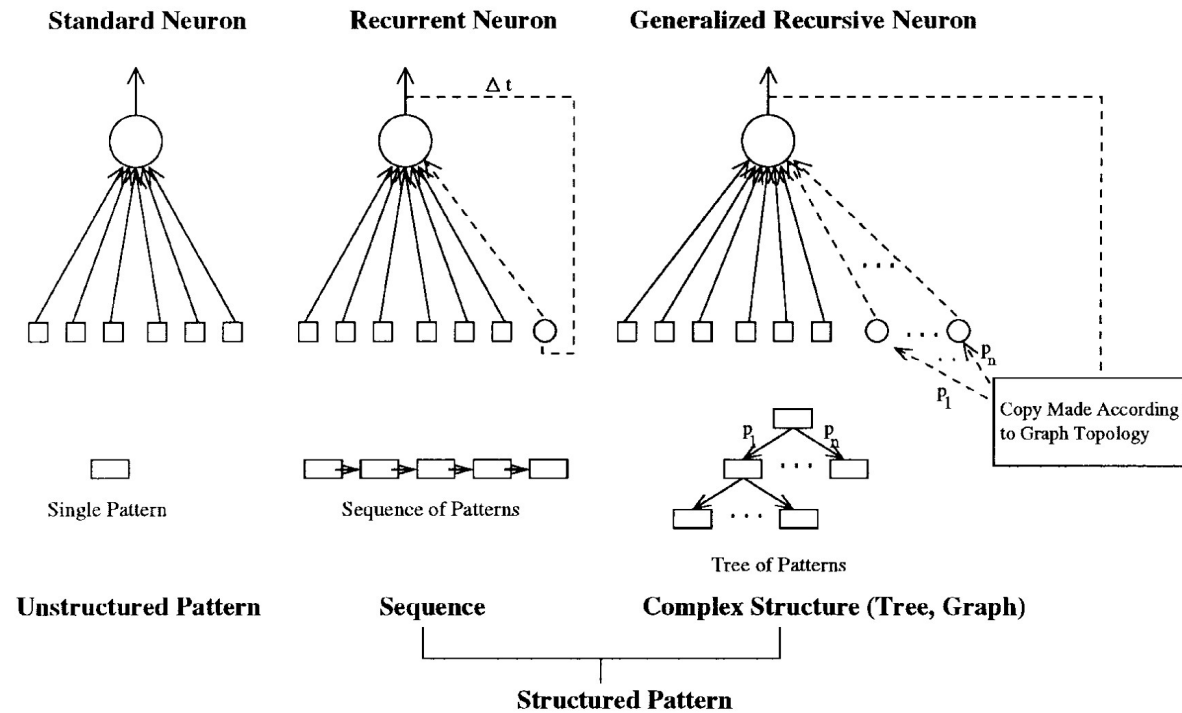


# The Model of Graph Neural Networks

---



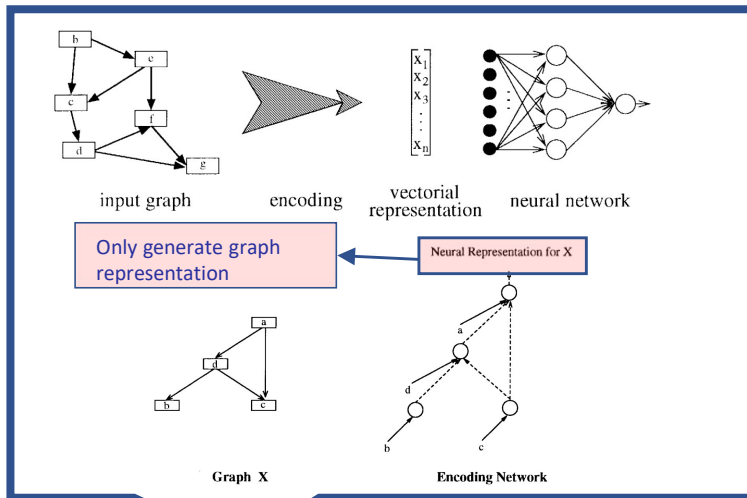
# GNN 1.0: Understanding GNN as RNN



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

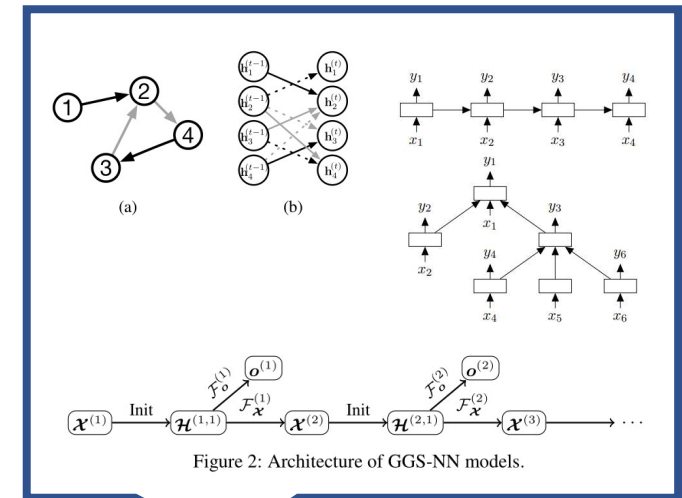
Sperduti, Alessandro, and Antonina Starita. 1997

# GNN 1.0: Understanding GNN as RNN



## From 2000 to 2010

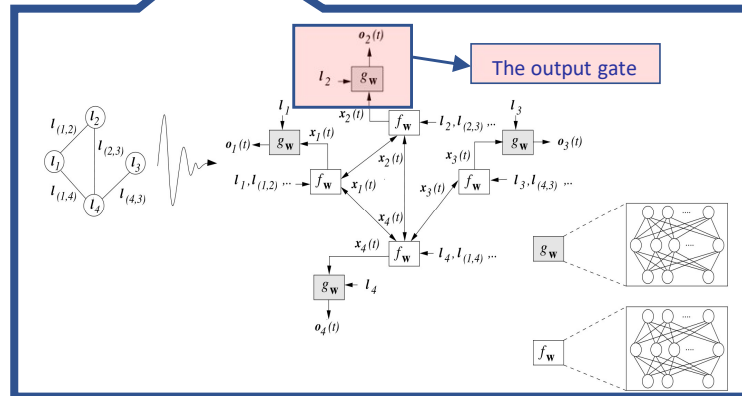
Gori et.al (IJCNN 05) and Scarselli et.al (TNN 08) add the **output gate** for each node to generate the node representation in graphs. This model is called GraphRNN.



## 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.

Sperduti, Alessandro, and Antonina Starita. 1997

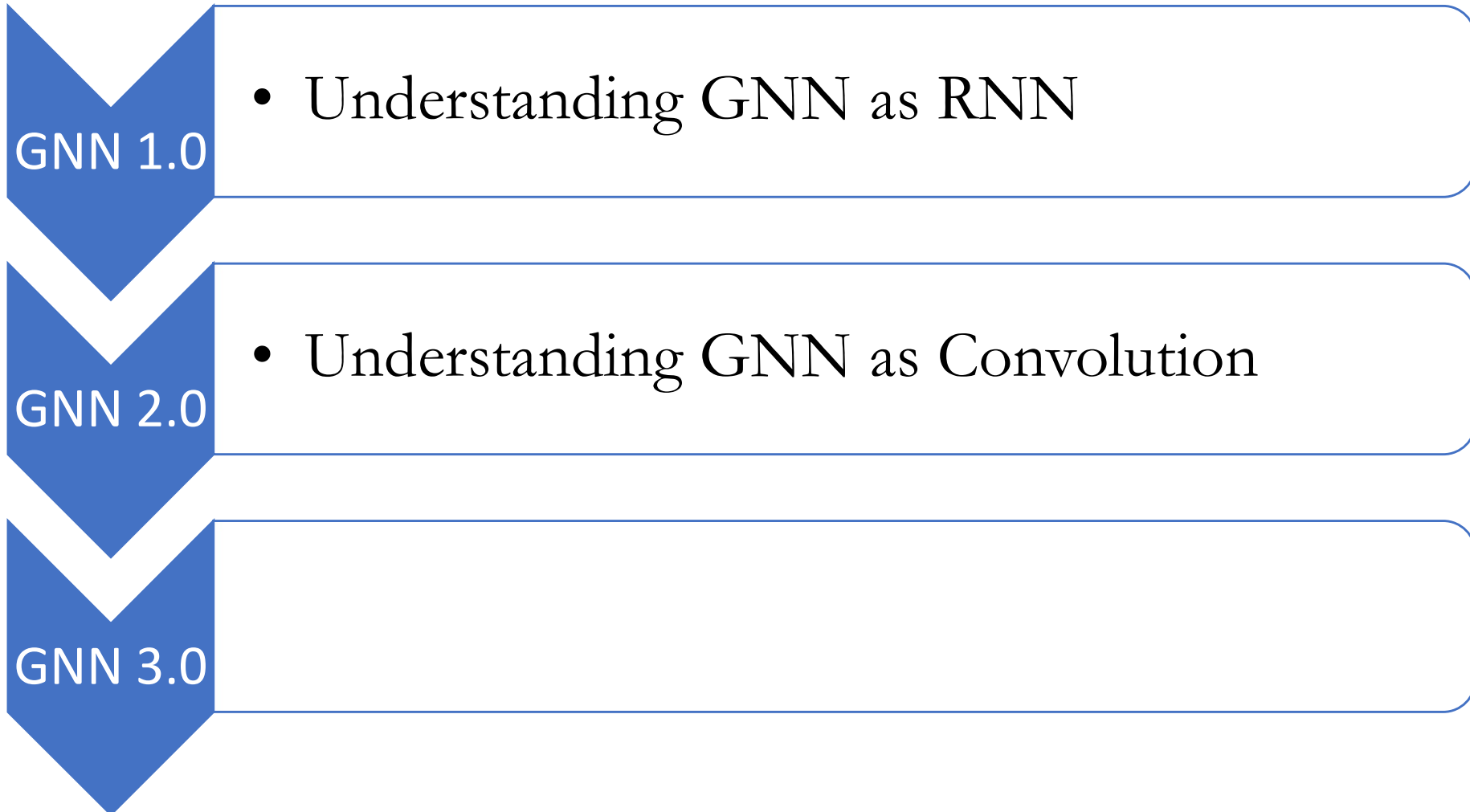
Gori, Marco, Gabriele Monfardini, and Franco Scarselli. 2005

Scarselli, Franco, et al. 2008

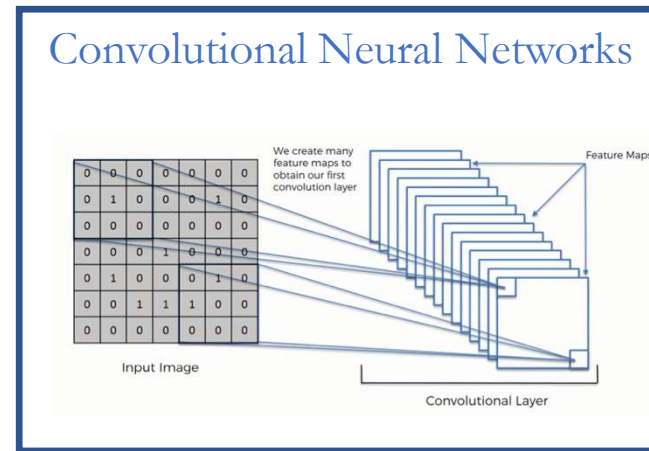
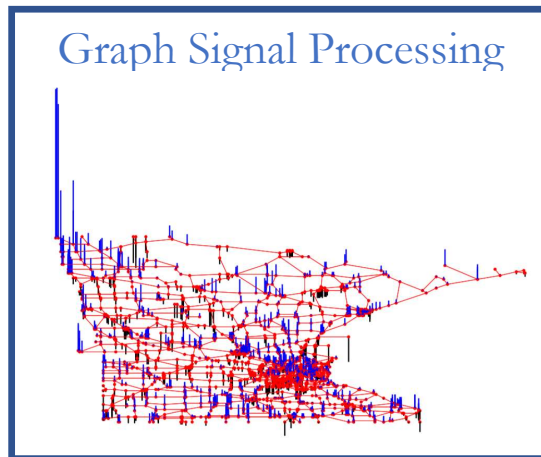
Li, Yujia, et.al. 2015, Tai, Kai Sheng et.al., 2015

# The Brief History of Graph Neural Networks

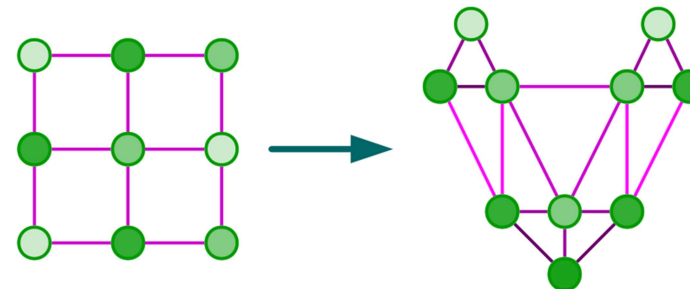
---



# 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]

Input graph signals e.g. bags of words → Feature extraction Convolutional layers → Classification Fully connected layers → Output signals e.g. labels

Graph signal filtering  
1. Convolution  
2. Non-linear activation

Graph coarsening  
3. Sub-sampling  
4. Pooling

$0 = \lambda_1 < \lambda_2 < \dots < \lambda_{M-1}$

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

[1] Bruna, Joan, et al. 2014  
 [2] Defferrard, Michaël, et al. 2016  
 [3] Niepert, Mathias, et al. 2016  
 [4] Kipf, Thomas N., and Max Welling. 2017

### Graph Convolutional Network (ICLR 2017)

input layer → hidden layers → output layer

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

- 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.**

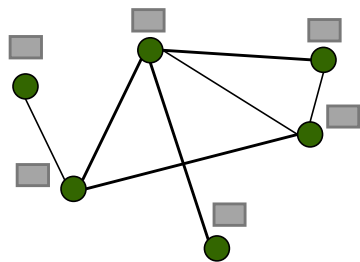
### PATCHY-SAN (ICML 2016)

- Neighborhood sampling to construct receptive field.

# Graph Signal Processing

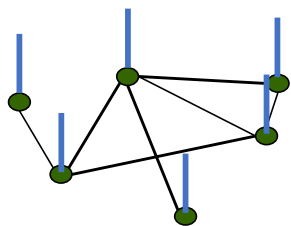
Graph signal:

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

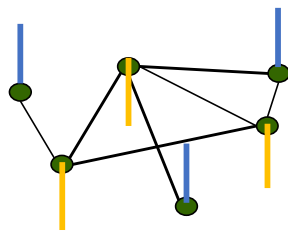


“Frequency” or “Smoothness” of the signal  $h$

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



Low frequency graph signal



High frequency graph signal

Eigen decomposition of graph Laplacian

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

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

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:

$$\mathbf{u}_i^T L \mathbf{u}_i = \mathbf{u}_i^T \lambda_i \mathbf{u}_i = \lambda_i$$

# Graph Convolution: Spectral domain $\rightarrow$ Spatial domain

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

$$y = x *_G \hat{g} = U \begin{pmatrix} \hat{g}(\lambda_1) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \hat{g}(\lambda_n) \end{pmatrix} U^T x = U \hat{g}(\Lambda) U^T x =: \hat{g}(L) x$$

**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  $\lambda_{max} = 2$ ,  $\theta := \theta_0 = -\theta_1$

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

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)

Spectral



Spatial

# 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

---

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

Lanczos Network [3]

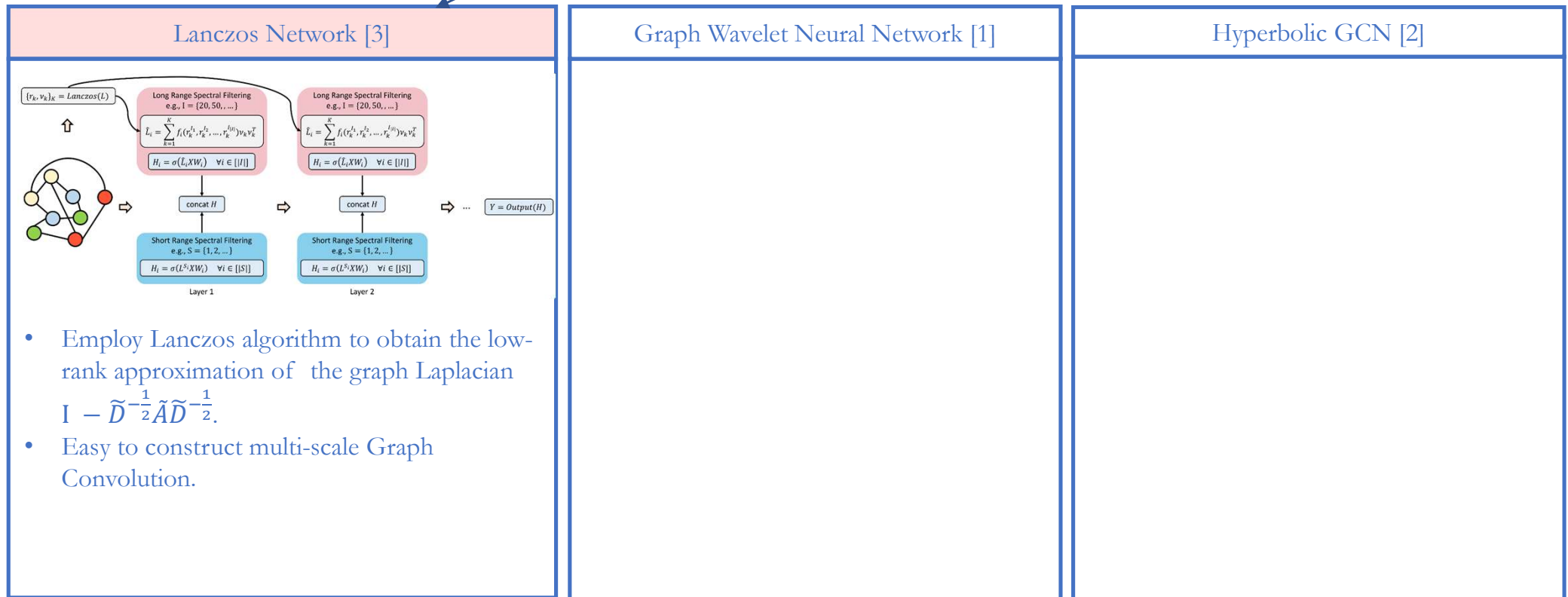
Graph Wavelet Neural Network [1]

Hyperbolic GCN [2]

[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

# GNN 3.0: Variants of Convolutions

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



[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

# 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]

$[v_k, v_k]_k = \text{Lanczos}(L)$   
 Long Range Spectral Filtering e.g.,  $l = \{20, 50, \dots\}$   
 $L_l = \sum_{k=1}^K f_l(v_k^l, v_k^l, \dots, v_k^{(l)}) v_k v_k^T$   
 $H_l = \sigma(L, XW) \quad \forall l \in [l]$   
 Short Range Spectral Filtering e.g.,  $S = \{1, 2, \dots\}$   
 $H_s = \sigma(L^S, XW) \quad \forall s \in [S]$   
 Layer 1      Layer 2  
 concat  $H$       concat  $H$        $Y = \text{Output}(H)$

- 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]

Wavelet Basis, scaling = 3      Wavelet Basis, scaling = 5  
 (a)      (b)

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]

[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

# 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]

$[v_k, v_k]_k = \text{Lanczos}(L)$   
 Long Range Spectral Filtering e.g.,  $l = \{20, 50, \dots\}$   
 $L_l = \sum_{k=1}^K f_l(v_k^1, v_k^2, \dots, v_k^{(l)}) v_k v_k^T$   
 $H_l = \sigma(L_l X W) \quad \forall l \in [L]$   
 Short Range Spectral Filtering e.g.,  $S = \{1, 2, \dots\}$   
 $H_s = \sigma(L^S X W) \quad \forall s \in [S]$

Layer 1:  $\text{concat } H \rightarrow \text{Short Range Spectral Filtering} \rightarrow H_s$   
 Layer 2:  $\text{concat } H \rightarrow \text{Short Range Spectral Filtering} \rightarrow H_s$   
 $Y = \text{Output}(H)$

- 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]

Wavelet Basis, scaling = 3 (a)  
 Wavelet Basis, scaling = 5 (b)

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]

$\mathbb{H}^{d,K}$   
 $\log_{x_j^H}^K$   
 $\mathcal{T}_{x_j^H} \mathbb{H}^{d,K}$   
 $\exp_{x_j^H}^K$   
 $\text{AGG}^{K_l}(x^H)_i$

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.
 
$$y_i^{(l+1),H} = \text{AGG}^{K_l}(h^{(l),H})_i$$

[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019



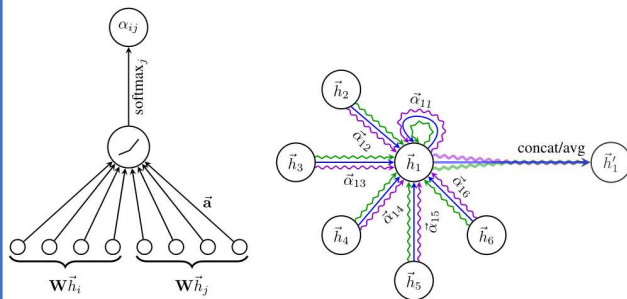
# GNN 3.0: GNN with Attention

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

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

Fixed during training

Graph Attention Network [1]



Replace the fixed aggregation weight  $\mathbf{a}_{ij}$  to the learnable self-attention.

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

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

Gated Attention Networks [2]

Spectral Graph Attention Network [3]

[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

# GNN 3.0: GNN with Attention

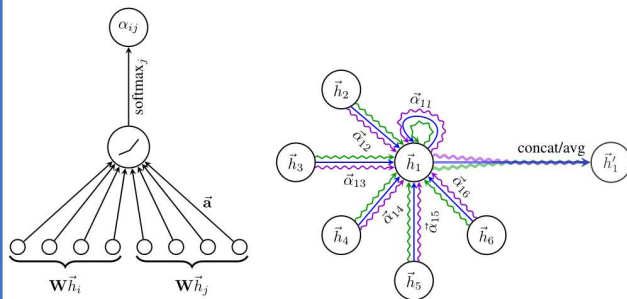
The original form:

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

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

Fixed during training

Graph Attention Network [1]

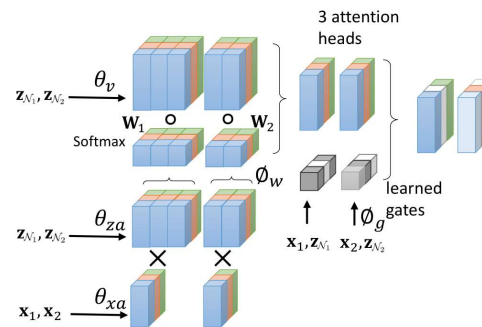


Replace the fixed aggregation weight  $a_{ij}$  to the learnable self-attention.

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

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

Gated Attention Networks [2]



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

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

$K$  is the number of heads.

Spectral Graph Attention Network [3]

[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

# GNN 3.0: GNN with Attention

The original form:

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

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

Fixed during training

### Graph Attention Network [1]

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

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

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

### Gated Attention Networks [2]

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

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

$K$  is the number of heads.

### Spectral Graph Attention Network [3]

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

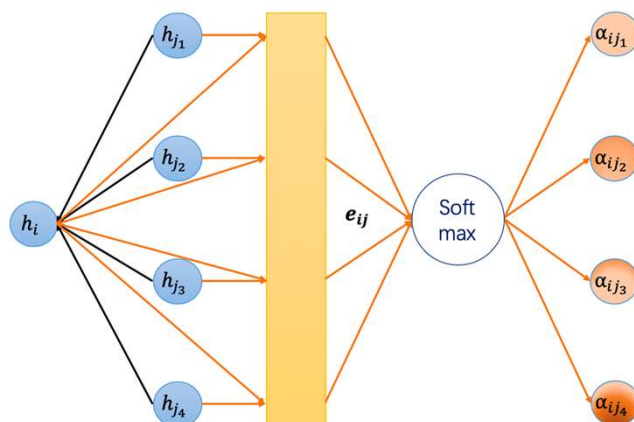
$$\mathbf{H}^{(l+1)} = \sigma(\text{AGG}(\mathbf{B}_L \alpha_L \mathbf{B}_L \mathbf{H}^{(l)}, \mathbf{B}_H \alpha_H \mathbf{B}_H \mathbf{H}^{(l)}) W^{(l)})$$

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

[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

# Graph Attention Network in Detail

## Single head attention



$$h_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} \alpha_{ij} W h_j^{(l)})$$

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

$$e_{ij} = \text{LeakyReLU}(\mathbf{a}^T (W h_i || W h_j))$$

## 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:

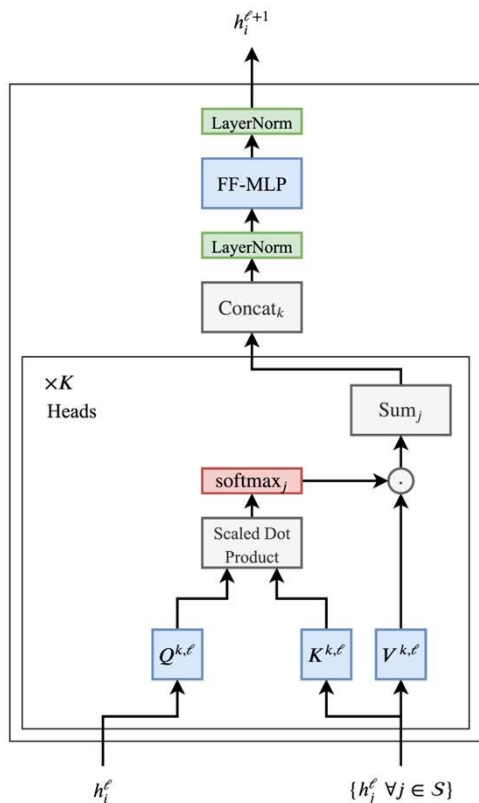
- Concatenation
- Average



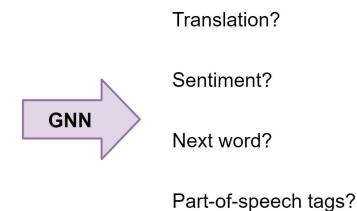
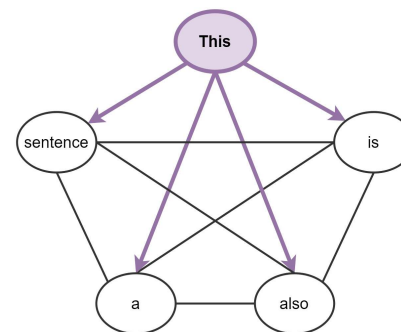
Attention weights learnt for the Cora dataset

# Transformers as GNNs with Multi-head Attention

One layer of the multi-head QKV attention



Transformer takes input sequence as a complete graph.



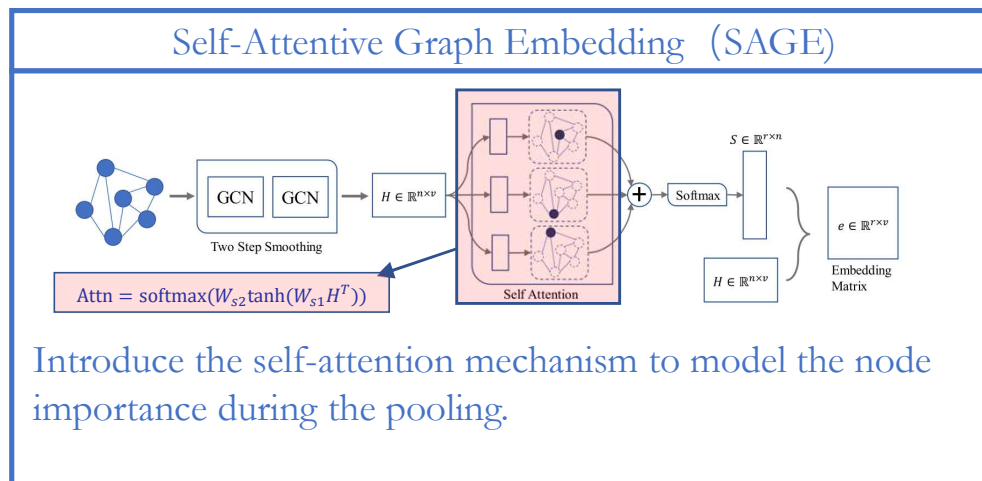
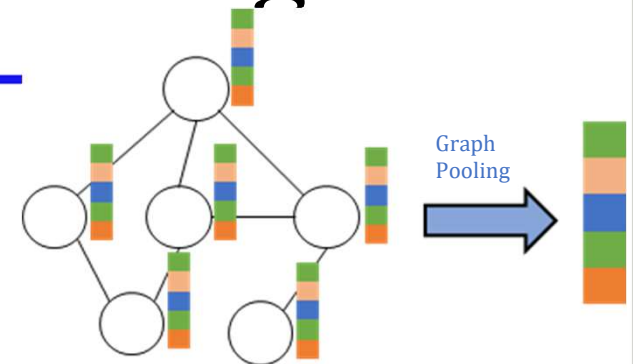
- ❑ Transformers can be viewed as GNNs with multi-head attention as the neighborhood aggregation function
- ❑ Transformers for NLP tasks treat the entire sequence as
- ❑ the neighborhood

Chaitanya Joshi. Transformers are graph neural networks, 2020.  
[https:// graphdeeplearning.github.io/post/transformers-are-gnns/](https://graphdeeplearning.github.io/post/transformers-are-gnns/)

# GNN 3.0: GNN with Graph Pooling

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

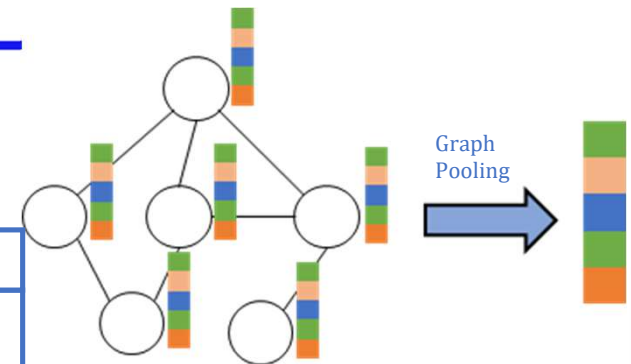
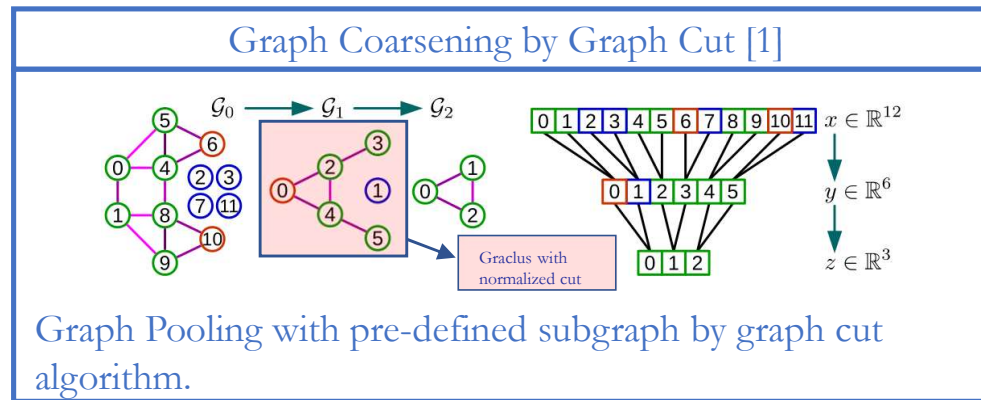
- The most straightforward way: Max/Mean Pooling
- SAGE: Attentive Pooling



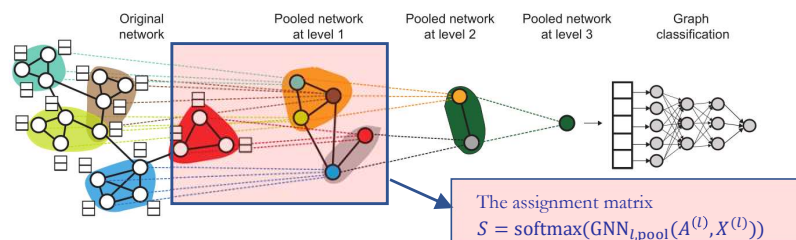
Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

# GNN 3.0: GNN with Graph Pooling

## Hierarchical Pooling

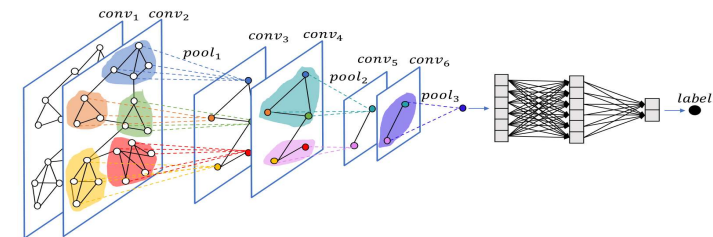


## Differentiable Graph Pooling (DIFFPOOL)[2]



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.

[1] Defferrard, Michaël, et al. 2016 [2] Ying, Zhitao, et al. 2018 [3] Ma, Yao, et al. 2019

# GNN Implementation: Message Passing Framework

- Message Passing Framework:

- Step 1: Gather and transform the messages from neighbors:

$$m_i^{(l+1)} = \text{AGG} \left( \{M^{(l+1)}(h_i^{(l)}, h_j^{(l)}, e_{i,j}) \mid j \in N(v_i)\} \right)$$

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

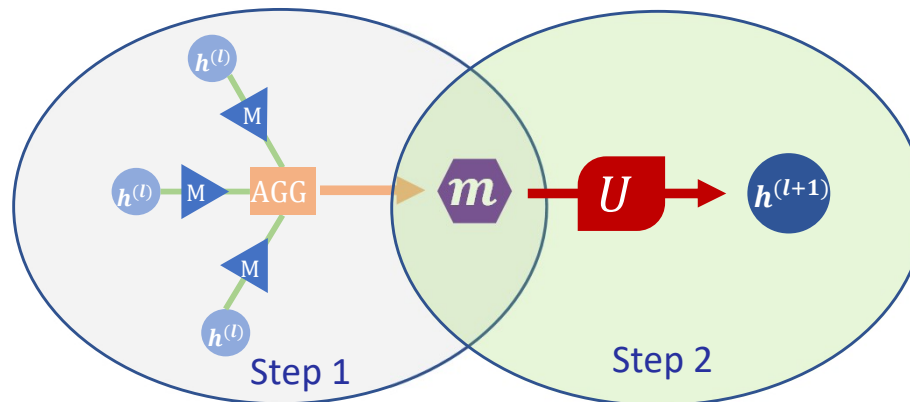
The neighborhood set of node. E.g. 1-hop neighbors.

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

$$h_i^{(l+1)} = U^{(l+1)}(h_i^{(l)}, m_i^{(l+1)})$$

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

The state update function.

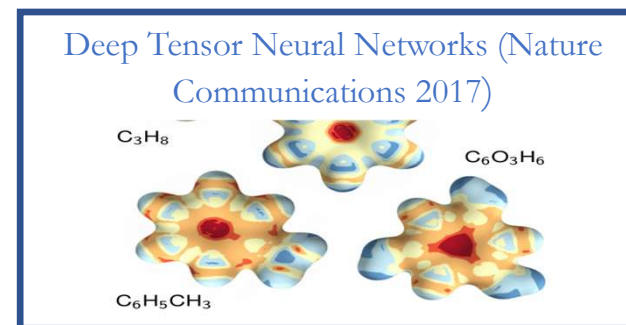
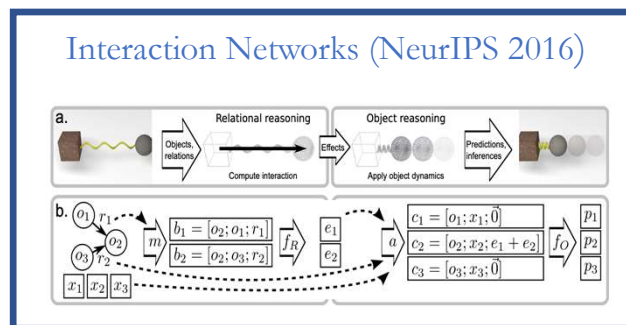
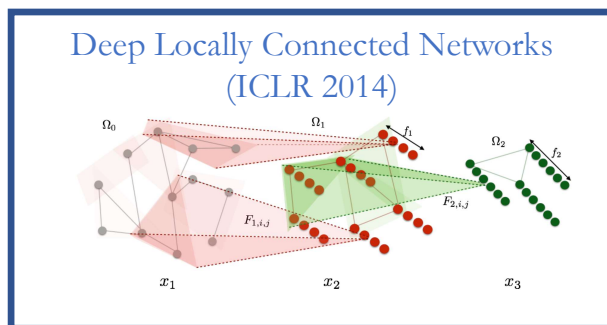
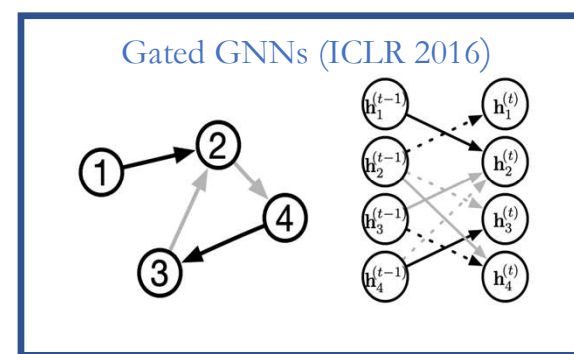
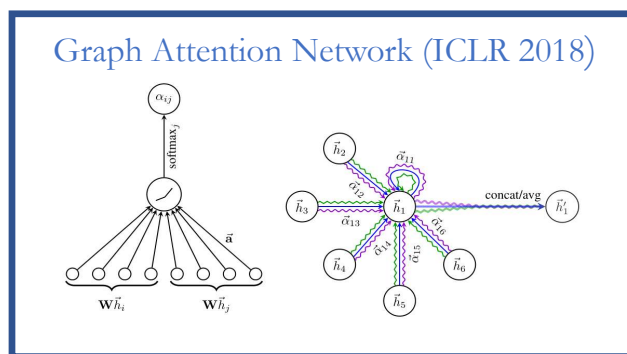
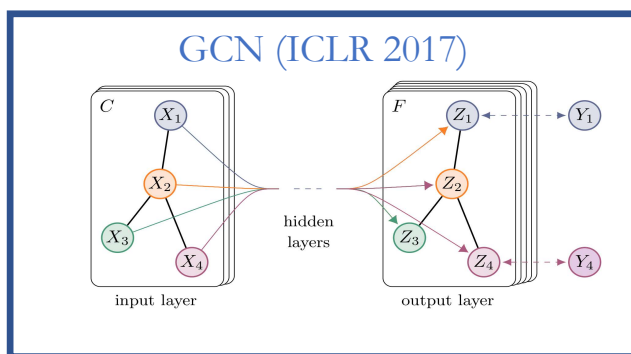


Gilmer, Justin, et al. "Neural Message Passing for Quantum Chemistry." ICML. 2017.



# Examples of Message Passing Realizations

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

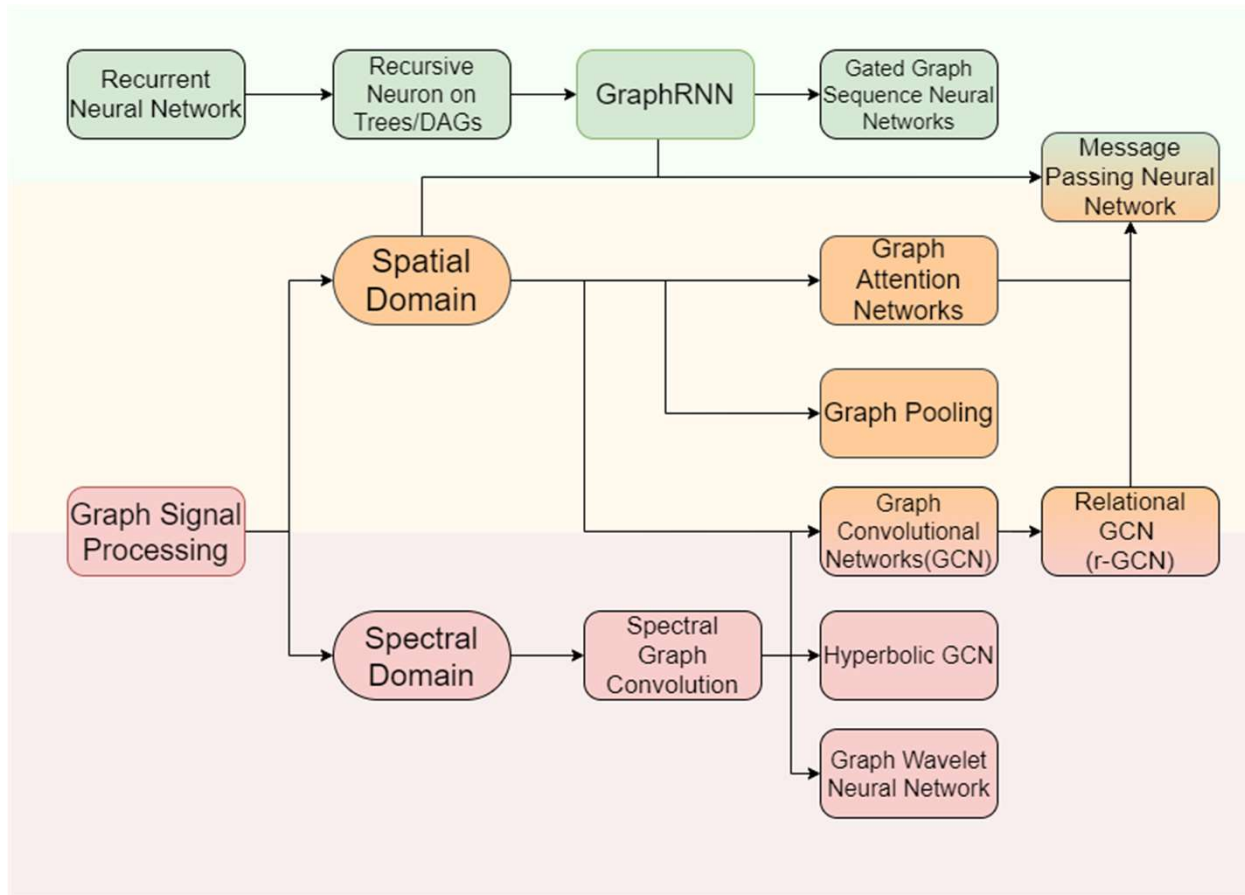


Wang, Minjie, et al. "Deep graph library: A graph-centric, highly-performant package for graph neural networks." *arXiv preprint arXiv:1909.01315* (2019).

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

# Summary

@change advanced topics



## Advanced topics

Training Deep GNNs

Scalability of GNNs

Robustness of GNNs

Self/Un-Supervised Learning of GNNs

....

---

# Training Deep GNNs

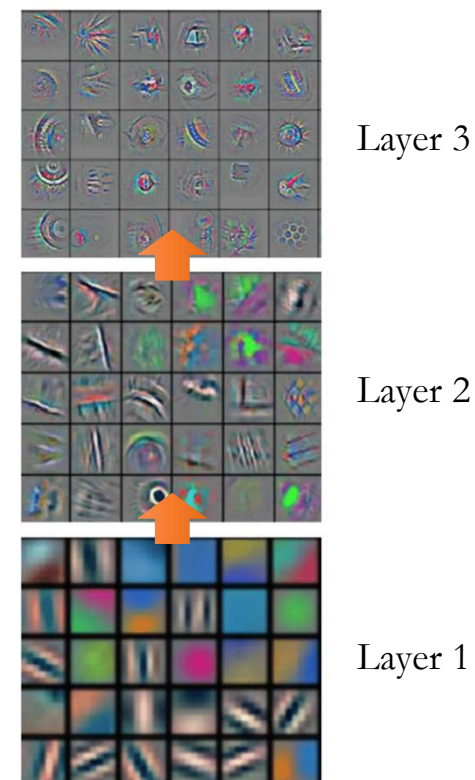
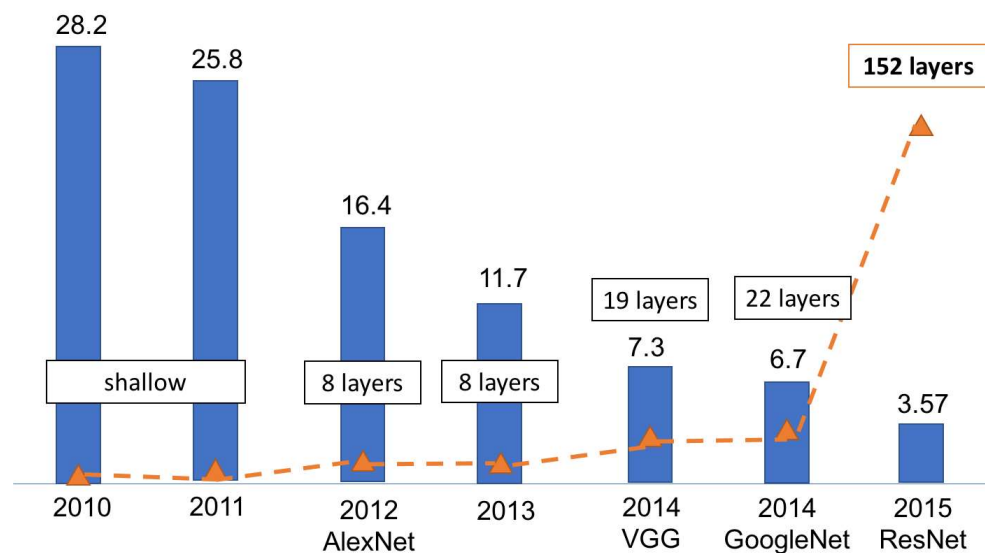
# Training Deep GNNs

---

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

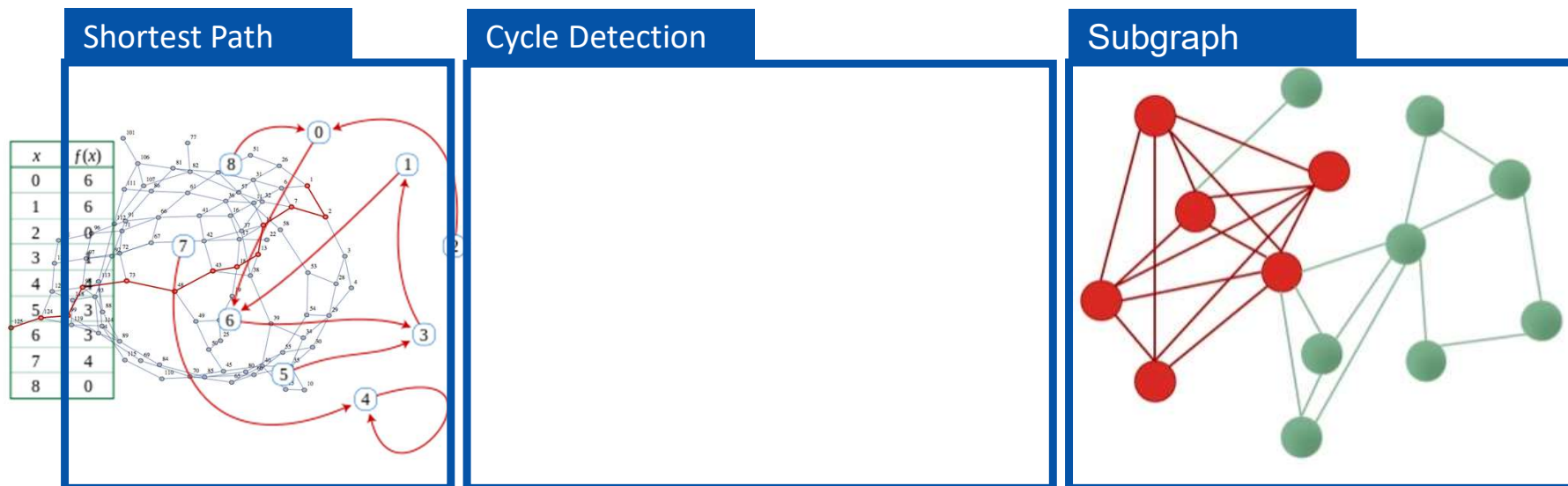
# The Power of Deep DNNs

- 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

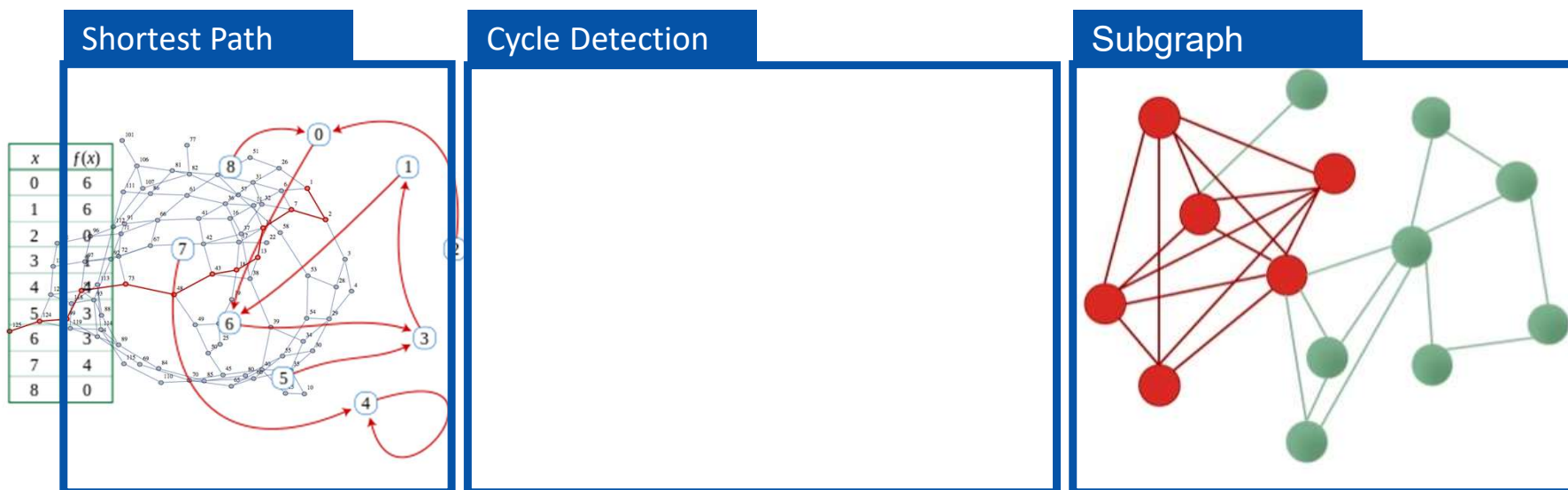


Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." *International Conference on Learning Representations*. 2020.

# 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



Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." *International Conference on Learning Representations*. 2020.

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

(Loukas, ICLR'20)

Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." *International Conference on Learning Representations*. 2020.



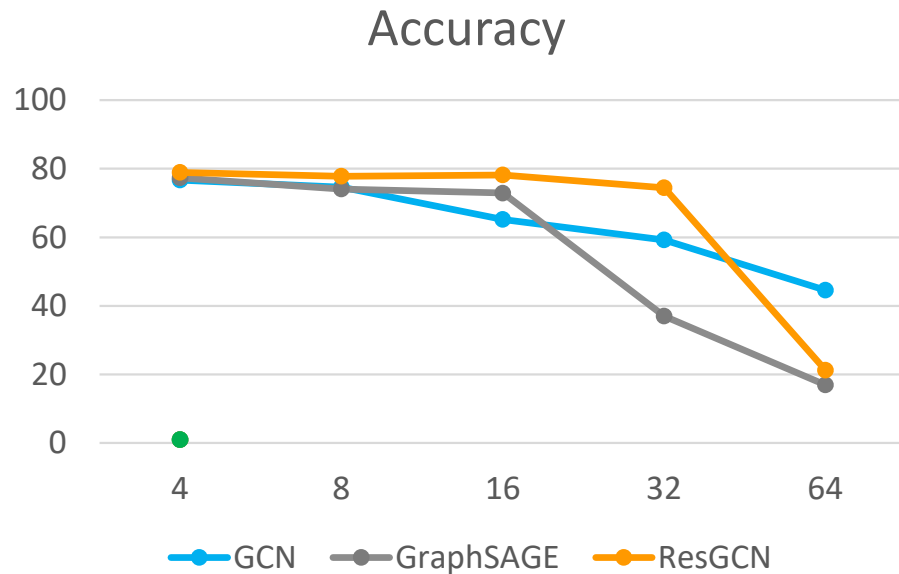
# Training Deep GNNs

---

- 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
<b>GCN</b>	<b>76.7</b>	<b>65.2</b>	<b>44.6</b>
<b>GraphSAGE</b>	<b>77.3</b>	<b>72.9</b>	<b>16.9</b>
<b>ResGCN</b>	<b>78.9</b>	<b>78.2</b>	<b>21.2</b>

(Rong et al, ICLR'20)

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 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?

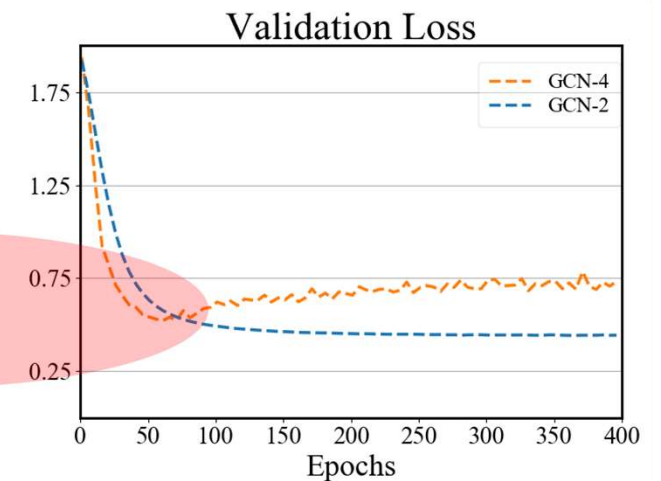
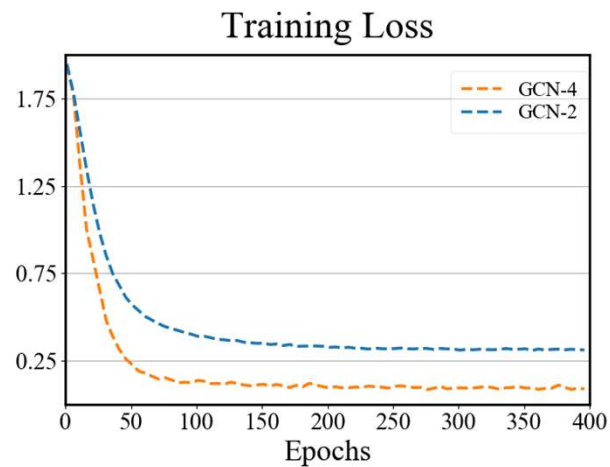
# Training deep GNNs

---

- 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

$$O(dh^2)$$

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

# Training deep GNNs

---

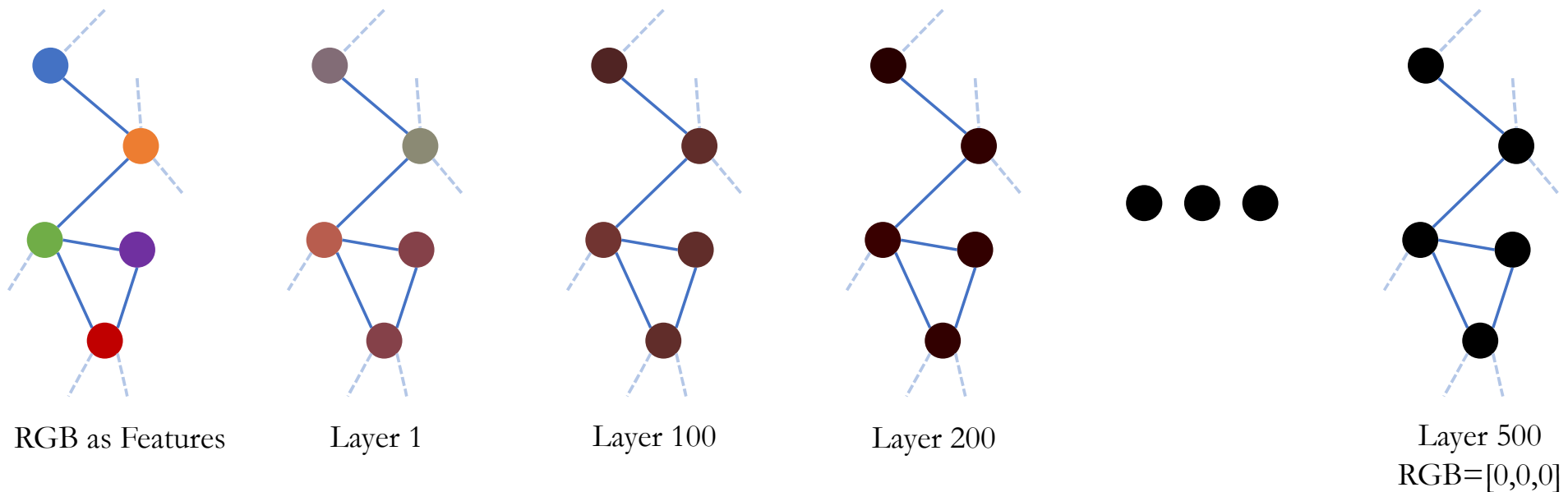
- Why do we need deep GNNs?
- Can GNNs simply go deep?
- What impedes GNNs to go deep?
  - Overfitting (Common)
  - Training dynamics (Common)**
  - Over-smoothing (Graph Specific)
- How to make GNNs deep?

# Training dynamics

$l$ -layers gradient

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

The gradients vanish as the model go deep because  $s_{1..l} \lambda_{m+1} < 1$



# 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?

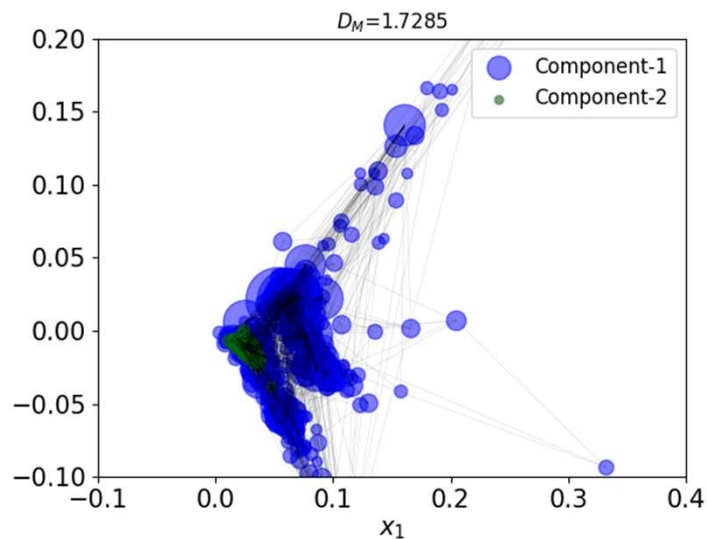


# Over-Smoothing

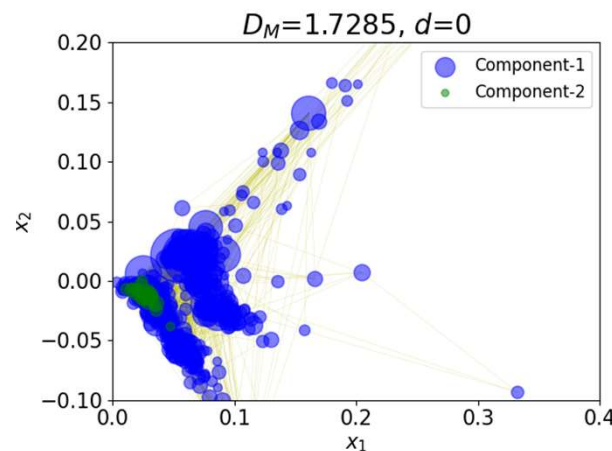
▶ GNNs suffers from over-smoothing

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

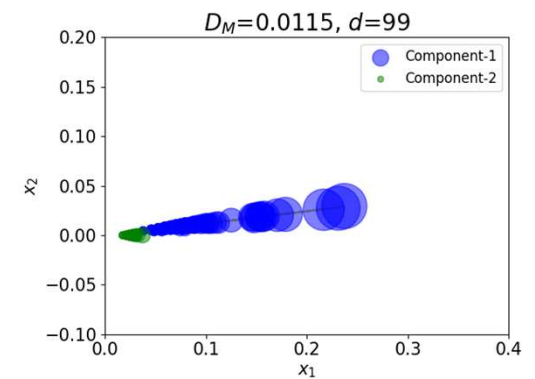
(Li et al. AAAI'18; Chen et al. AAAI'20; Oono et al. ICLR'20; Rong et al. ICLR'20; Huang et al. arXiv'20)



(Huang et al. arXiv'20)



Initial

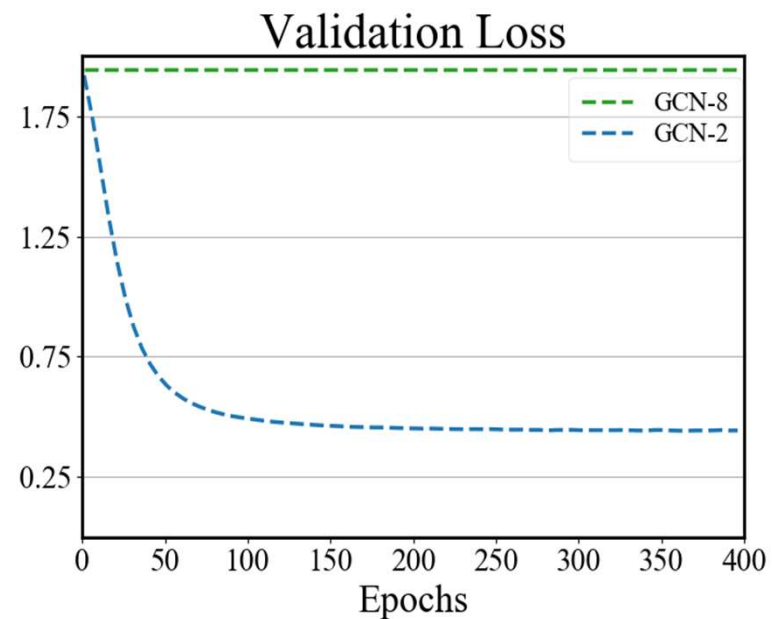
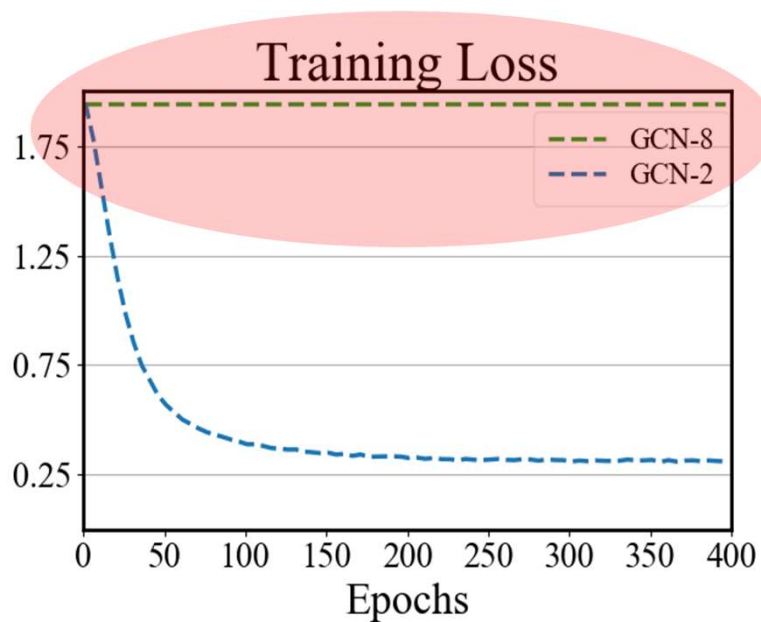


99 layers

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Over-Smoothing

Over-smoothing also impedes training.



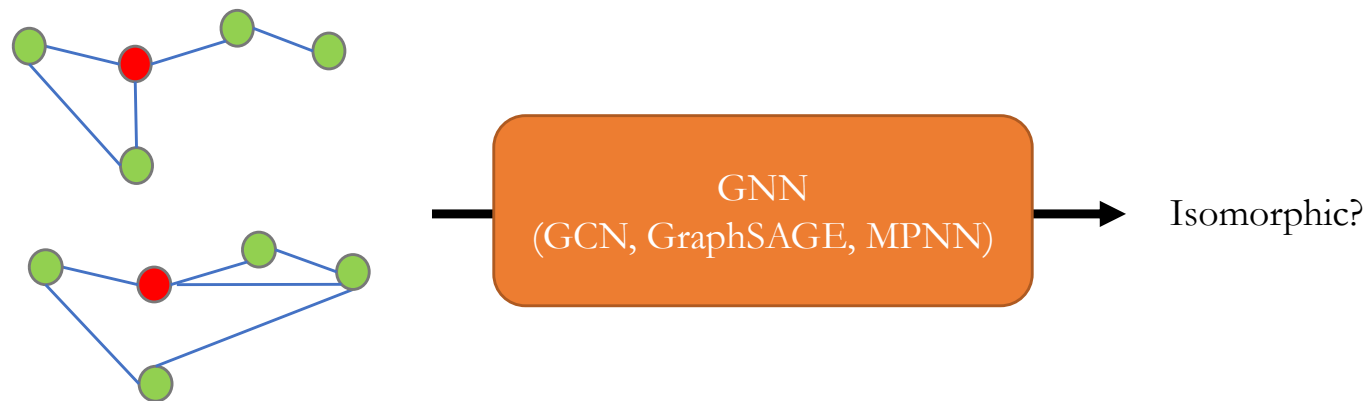
Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Over-Smoothing

---

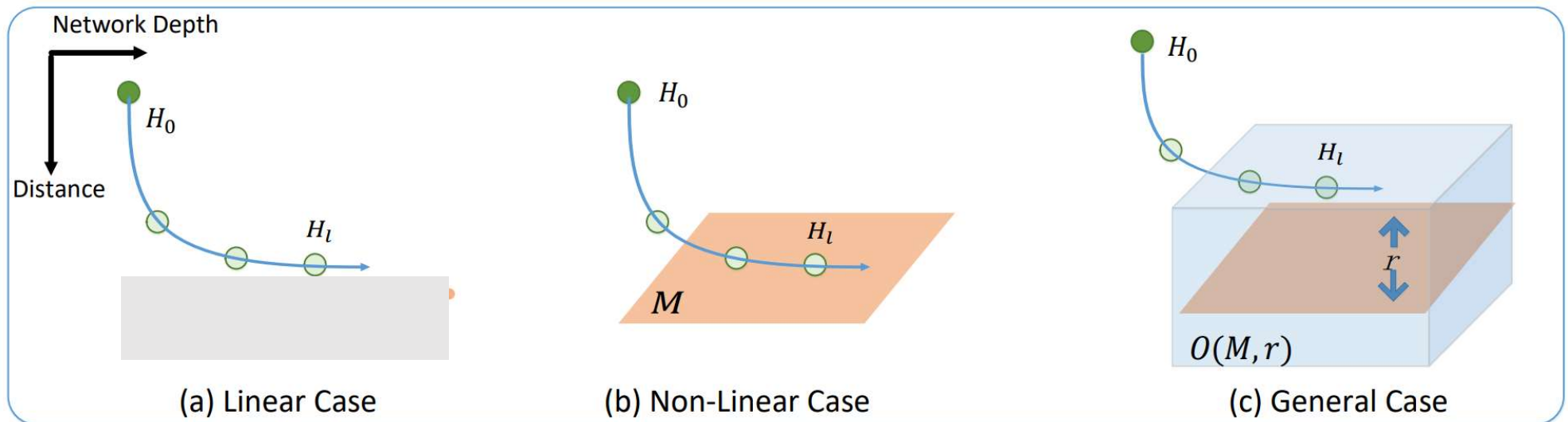
## Why GNN works?

- 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



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

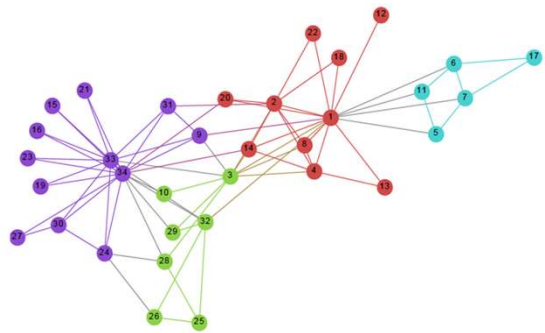
# Over-Smoothing of Linear GCN

- When GCNs fail?
- With linear activation

$l$ -step Random Walk

Probability of walking

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



Random Walks on Graph

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

Tang, Jian, et al. "Line: Large-scale information network embedding." In WWW, 2015.

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAAI/2018.

# 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

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." *AAAI*/2018.

# 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

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

Eigen decomposition

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

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAAI/2018.

# Over-Smoothing of Linear GCN

---

👉 Rewrite eigen decomposition

Eigen decomposition

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

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAAI/2018.



# Over-Smoothing of Linear GCN

---

Rewrite eigen decomposition

Eigen decomposition

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

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

Eigenvalues

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

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAAI/2018.

# Over-Smoothing of Linear GCN

Rewrite eigen decomposition

Eigen decomposition

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

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

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

1

1

0

0

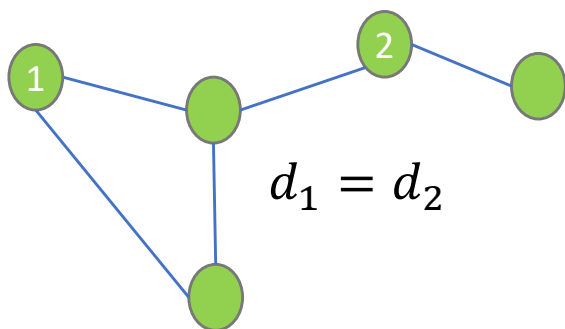
# Over-Smoothing of Linear GCN

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

## Convergence

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

Node 1 is **indistinguishable** with node 2



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

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

# Over-Smoothing of Non-Linear GCN

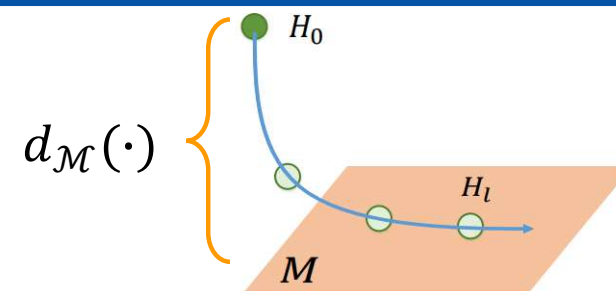
With ReLU activation

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

$\mathcal{M}$  subspace

**Definition 1** (subspace). Let  $\mathcal{M} := \{EC \mid C \in \mathbb{R}^{M \times C}\}$  be an  $M$ -dimensional subspace in  $\mathbb{R}^{N \times C}$ , where  $E \in \mathbb{R}^{N \times M}$  is orthogonal, i.e.  $E^T E = I_M$ , 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  $\tilde{A}$

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



Convergence

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

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *ICLR 2020*.

# Over-Smoothing of Non-Linear GCN

---

Convergence of  $\tilde{A}$

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

Convergence of  $W$

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



Convergence

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

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *ICLR 2020*.

# Over-Smoothing of Non-Linear GCN

Convergence of  $\tilde{A}$

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

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

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

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *ICLR 2020*.

# Over-Smoothing of GCNs with bias

---

- With ReLU and bias

GCNs with bias

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

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020



# Over-Smoothing of GCNs with bias

---

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

GCNs with bias

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

Convergence of bias

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

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# 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

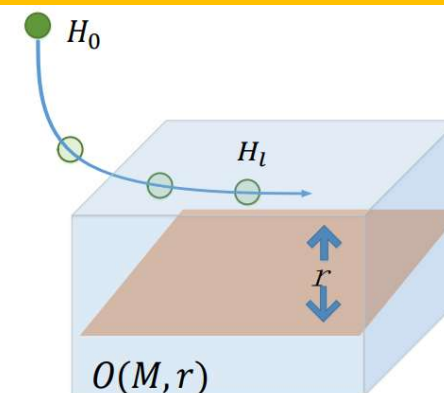
$$d_{\mathcal{M}}(H_{l+1}) \leq \lambda_{m+1}s d_{\mathcal{M}}(H_l) + d_{\mathcal{M}}(b)$$

GCN with bias

$H_l$  converges to a certain sub-cube:

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

(Huang et al. arXiv'20)



(c) General Case

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# 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

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Other methods to measure over-smoothing

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

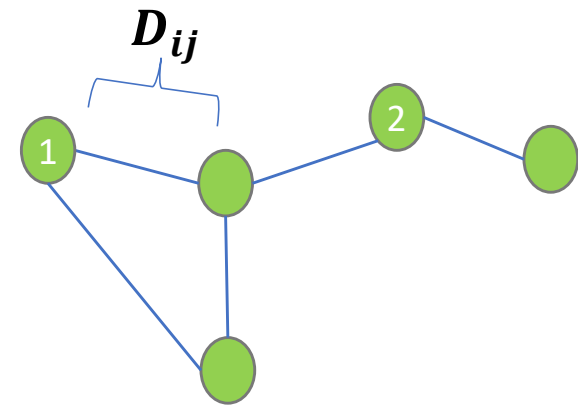
Pair Distance

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

The i-th row of the output

Mean Average Distance

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



- 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.

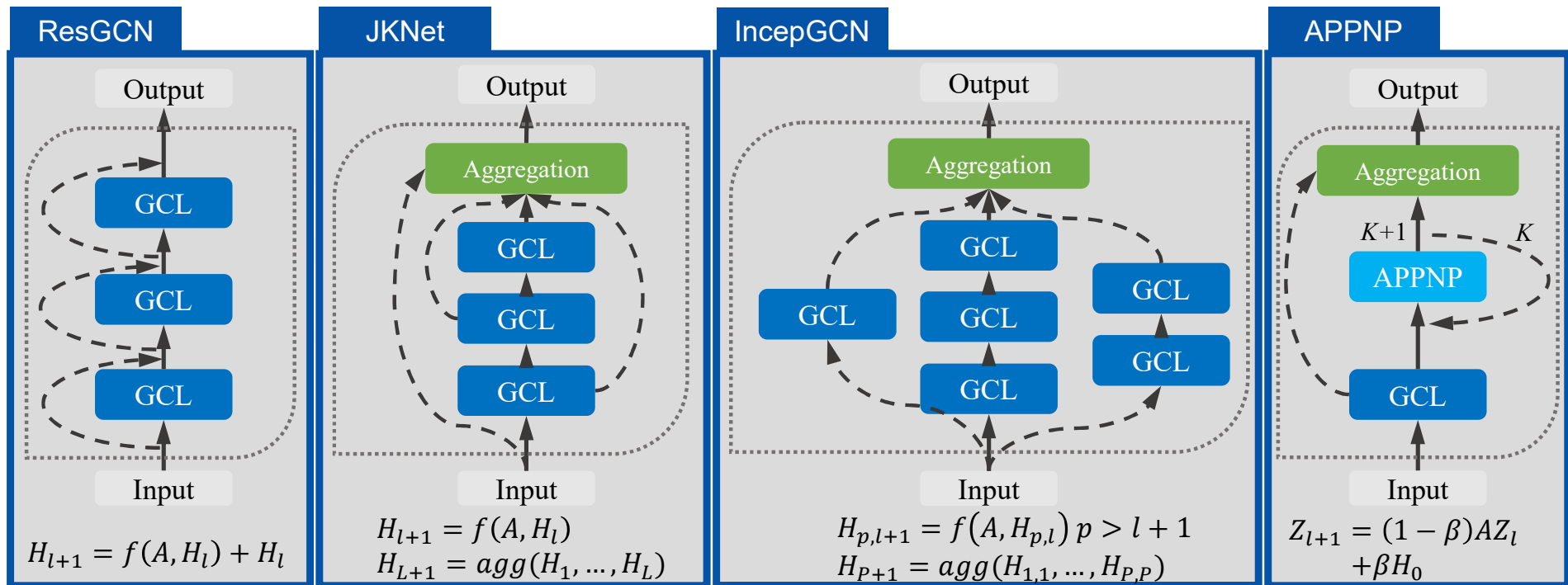
Zhou et al. "Towards Deeper Graph Neural Networks with Differentiable Group Normalization Kaixiong Zhou Texas A&M University zkxiong@tamu.edu Xiao Huang The Hong." NIPS 2020.

# 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

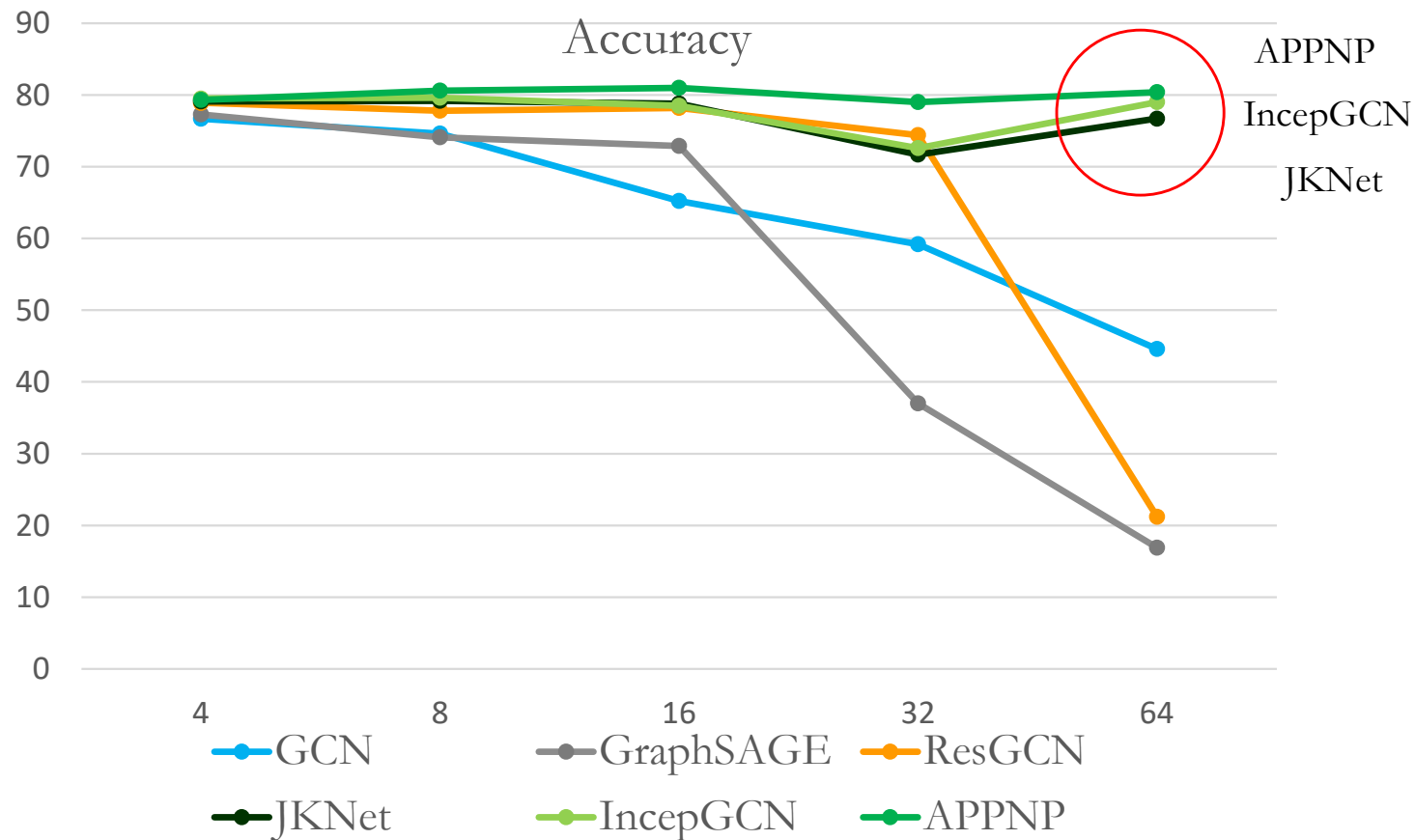
# Shortcuts in Structures



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

Rong, Yu, et al. "Droptedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

# Shortcuts in Structures



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

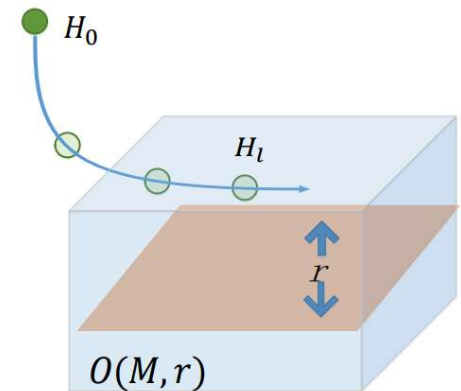
# Shortcuts in Structures

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

GCN with bias

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

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



(c) General Case

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020



# Shortcuts in Structures

## General Case

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

### Basic GCN

#### Generic GCN

$$v = s\lambda_{m+1}$$
$$r = 0$$

#### GCN with bias

$$v = s\lambda_{m+1}$$
$$r = \frac{d_{\mathcal{M}}(b)}{1 - v}$$

### Different Structures

#### ResGCN

$$v = s\lambda_{m+1} + \alpha$$
$$r = 0$$

#### APPNP

$$v = (1 - \beta) \lambda_{m+1}$$
$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - v}$$

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Shortcuts in Structures

## General Case

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

### Basic GCN

#### Generic GCN

$$v = s\lambda_{m+1}$$
$$r = 0$$

#### GCN with bias

$$v = s\lambda_{m+1} \quad r \text{ is nonzero}$$
$$r = \frac{d_{\mathcal{M}}(\mathbf{b})}{1 - v}$$

### Different Structures

#### ResGCN

$$v = s\lambda_{m+1} + \alpha$$
$$r = 0$$

#### APPNP

$$v = (1 - \beta) \lambda_{m+1}$$
$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - v}$$

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Shortcuts in Structures

## General Case

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

### Basic GCN

#### Generic GCN

$$v = s\lambda_{m+1}$$
$$r = 0$$

#### GCN with bias

$$v = s\lambda_{m+1}$$
$$r = \frac{d_{\mathcal{M}}(b)}{1 - v}$$

### Different Structures

#### ResGCN

$v$  is increased

$$v = s\lambda_{m+1} + \alpha$$
$$r = 0$$

#### APPNP

$$v = (1 - \beta) \lambda_{m+1}$$
$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - v}$$

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Shortcuts in Structures

## General Case

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

### Basic GCN

#### Generic GCN

$$v = s\lambda_{m+1}$$
$$r = 0$$

#### GCN with bias

$$v = s\lambda_{m+1}$$
$$r = \frac{d_{\mathcal{M}}(b)}{1 - v}$$

### Different Structures

#### ResGCN

$$v = s\lambda_{m+1} + \alpha$$
$$r = 0$$

#### APPNP

$$v = (1 - \beta) \lambda_{m+1}$$
$$r = \frac{\beta d_{\mathcal{M}}(H_0)}{1 - v}$$

$v$  is decreased  
 $r$  is nonzero

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Over-smoothing Layer

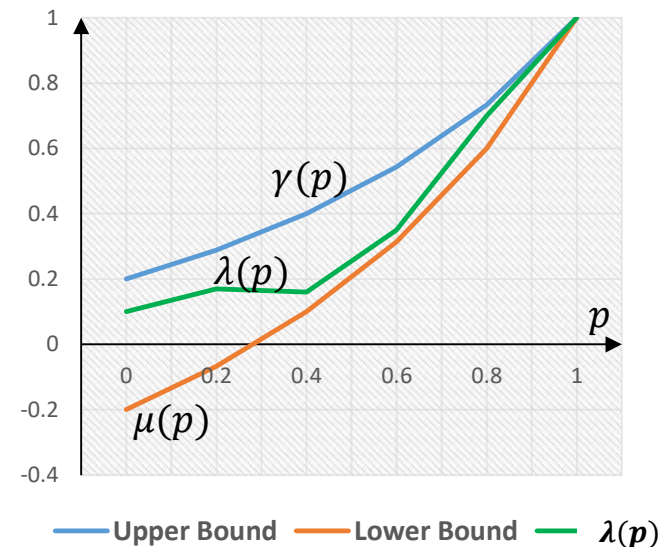
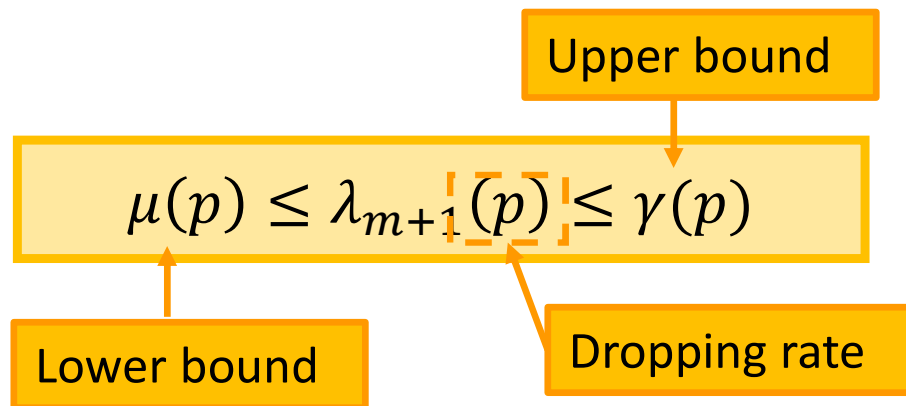
---

- 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.  $p$ ;

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# 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)

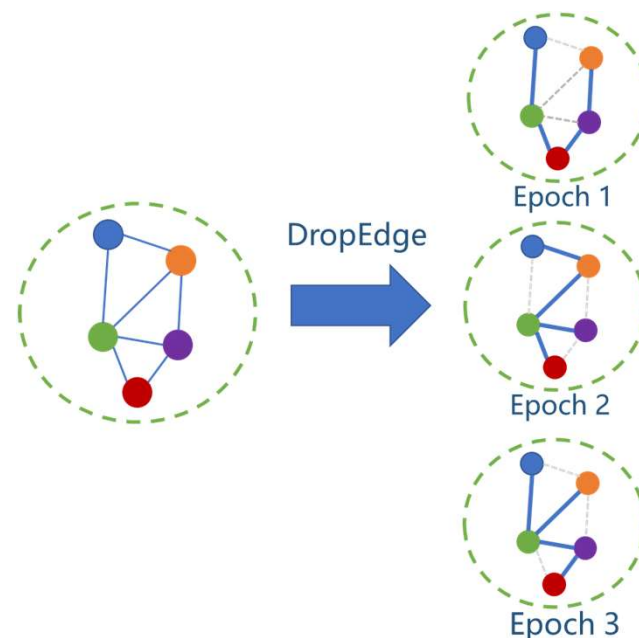
Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Alleviate Over-Smoothing by DropEdge

---

Besides, DropEdge can prevent over-fitting as well!



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020



# Alleviate Over-Smoothing by Adjacency Matrix

## DropEdge results

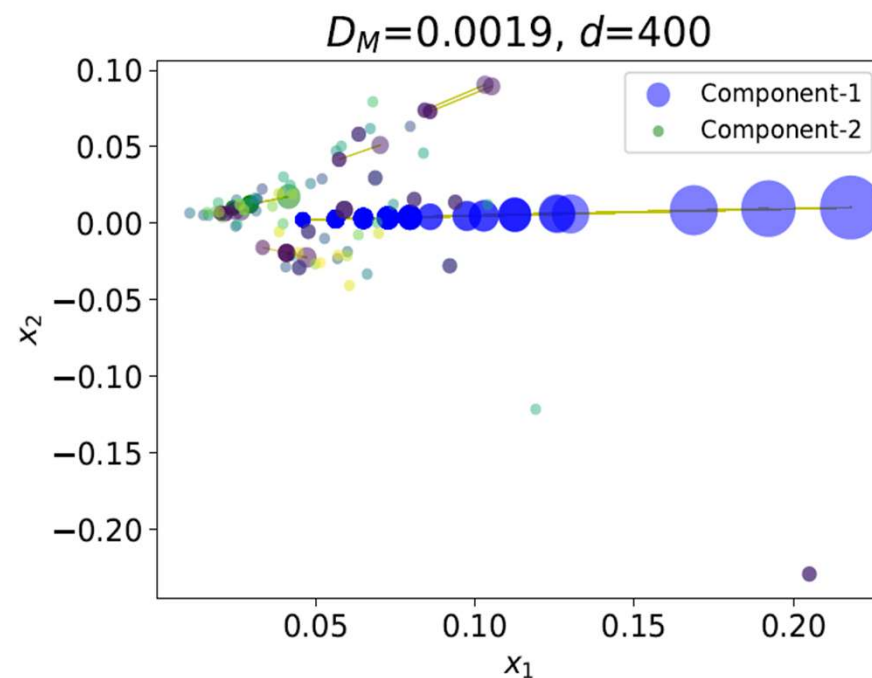
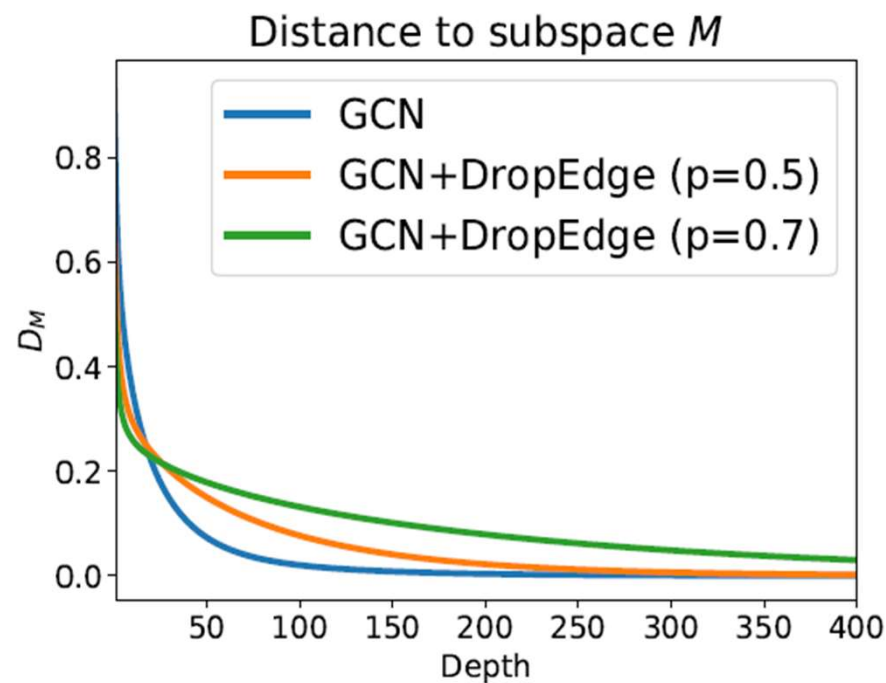
Citeseer	4 layers	DropEdges	16 layers	DropEdges	64 layers	DropEdges
GCN	76.7	79.2(+2.5)	65.2	76.8(+11.6)	44.6	45.6(+1.0)
ResGCN	78.9	78.8(-0.1)	78.2	79.4(+1.2)	21.2	75.3(+54.1)
JKNet	79.1	80.2(+1.1)	78.8	80.1(+1.3)	76.7	80.0(+3.3)
IncepGCN	79.5	79.9(+0.4)	78.5	80.2(+1.7)	79.0	79.9(+0.9)
GraphSAGE	77.3	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)

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Alleviate Over-Smoothing by Adjacency Matrix

## DropEdge results



Rong, Yu, et al. "Dropege: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

# Alleviate Over-Smoothing by Weights

---

Convergence speed

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

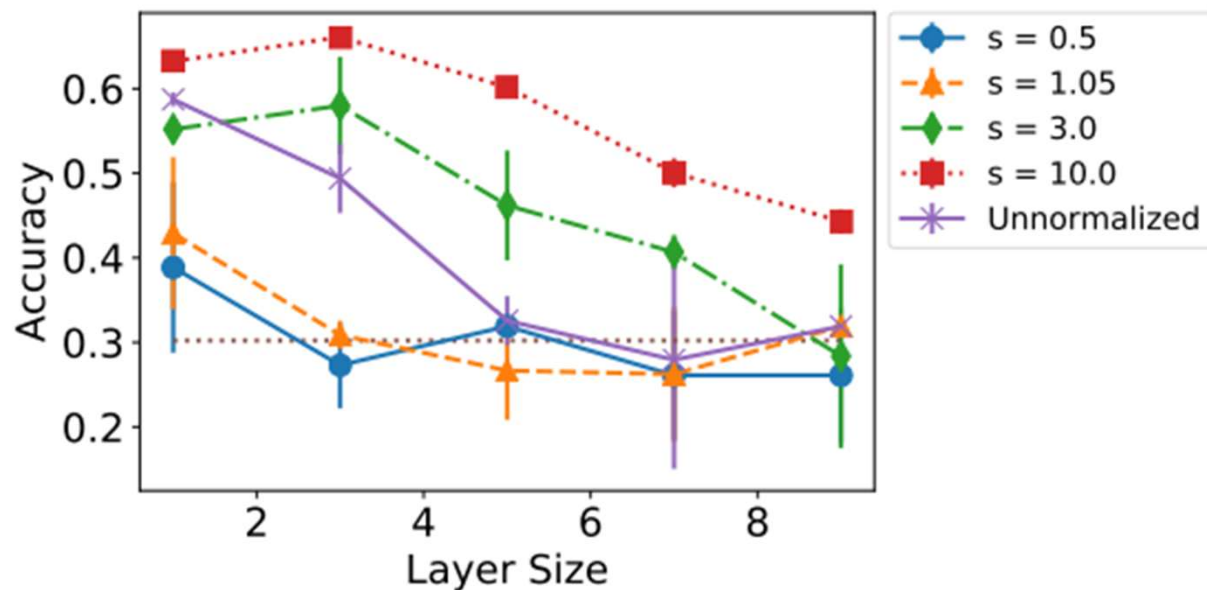
Weights

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

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *ICLR 2020*.

# Alleviate Over-Smoothing by Weights

Try different  $s$  as initial



Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." *ICLR 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

# Pair Norm: Center and Rescale

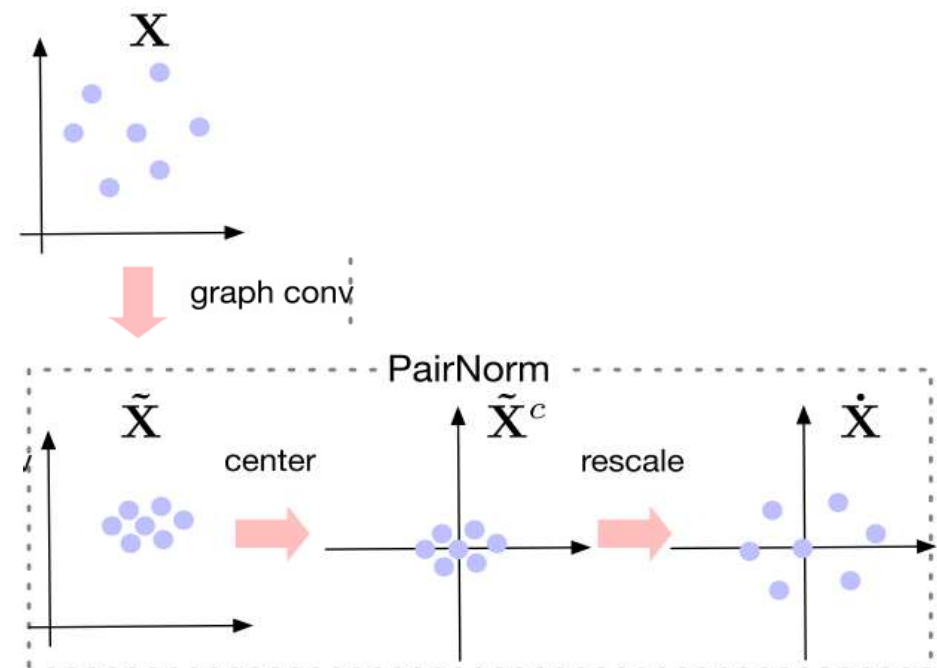
- PairNorm: Center and rescale (normalize) GCN outputs  $\tilde{X} := \text{GCN}(A, X)$  to keep the **total pairwise squared distance** *unchanged*

Center

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

Rescale

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



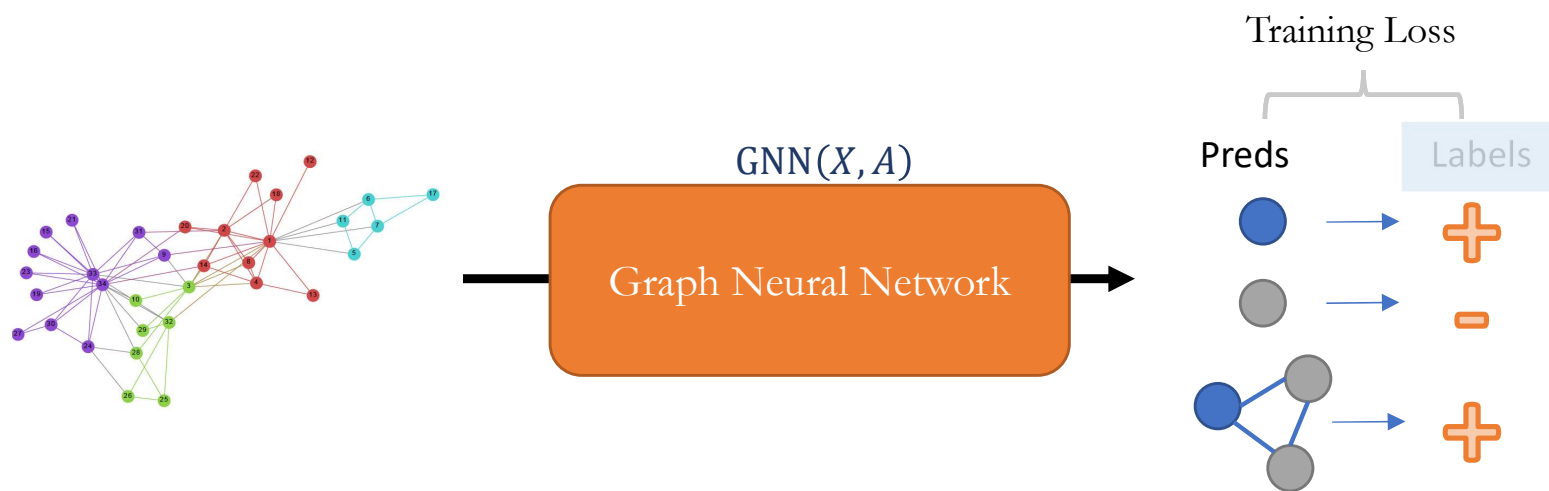
- See also GroupNorm (Zhou et al., 2020)

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

---

# Self/Un-Supervised Learning of GNNs

# What we discussed before are supervised

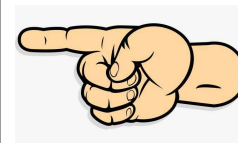


- **Labels are scarce**, e.g. molecular property
- **Training/Testing tasks are Non I.I.D.**



# Existing Self-Supervised GNNs

	Node-Classification	Link/Metapath Prediction	Graph-Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	<b>S<sup>2</sup>GRL</b> [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

---

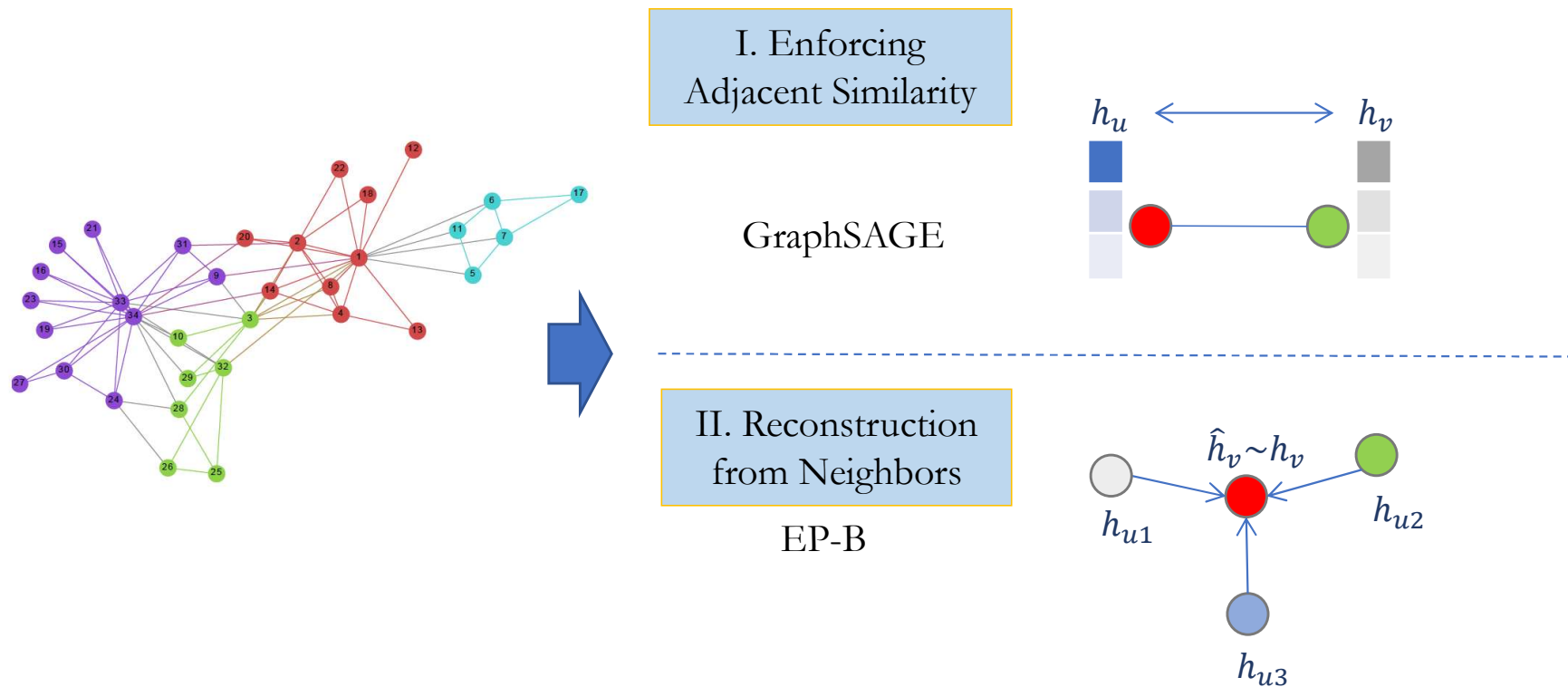
“In self-supervised learning, the system learns to predict part of its input from other parts of its input.” ---- by Yann Lecun

**Graphs are highly structured!**



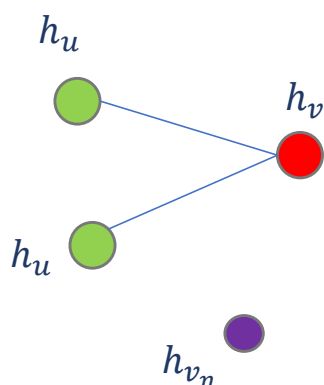
# Node Classification

- Two typical ways to formulate training loss



# I Enforcing Adjacent Similarity

- GraphSAGE (Hamilton et al. 2017)



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

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

Positive Samples

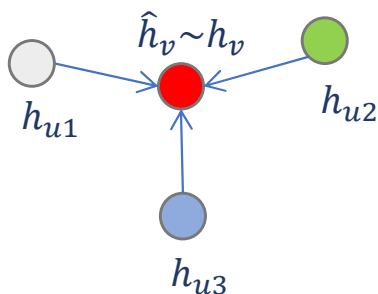
Negative Samples

$h_v$ : 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):



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

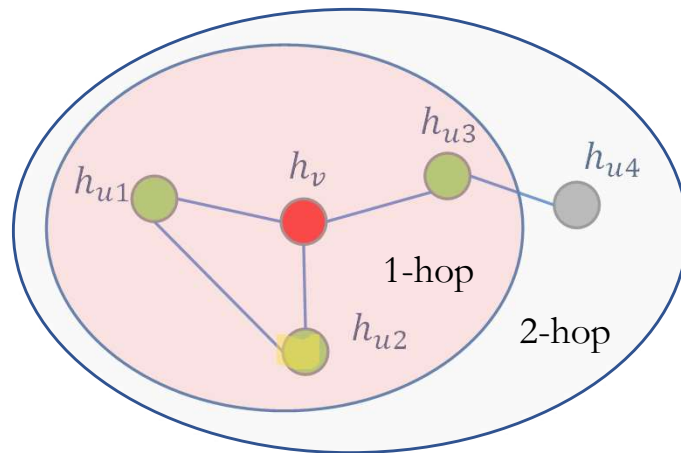
Positive Samples
Negative Samples

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

Durán & Niepert. Learning Graph Representations with Embedding Propagation.

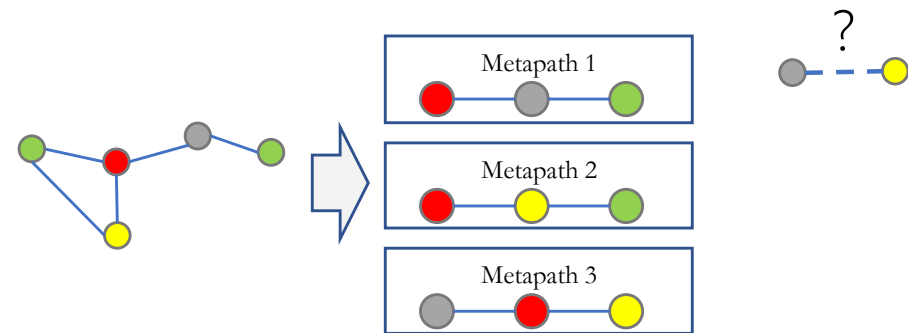
# Link/Metapath Prediction

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



Predict hop counts (K-hop connectivity)

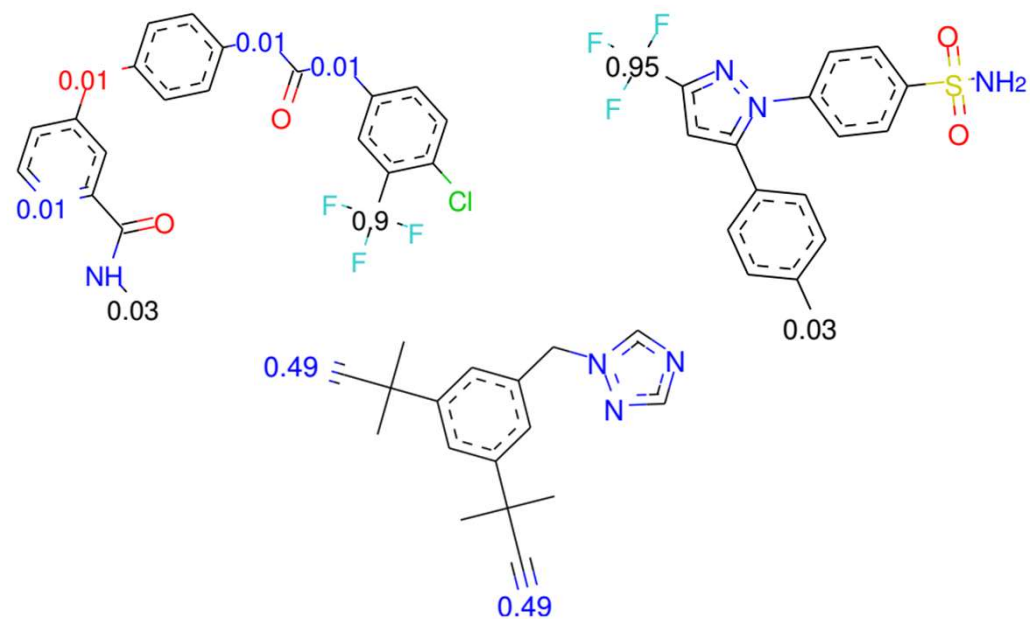
- SELAR(Hwang, Dasol, et al. 2020 )



Predict the type of meta path between two nodes

# How about graph classification/regression?

---

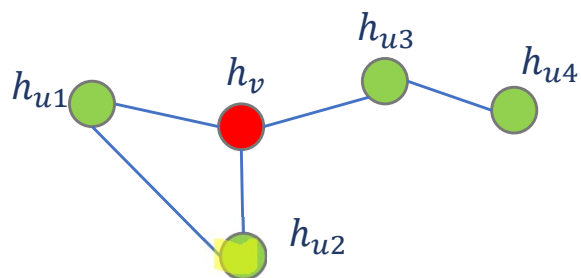


Toxicity?  
Solubility?  
...

# N-Gram Graph

- (Liu et al. 2019)

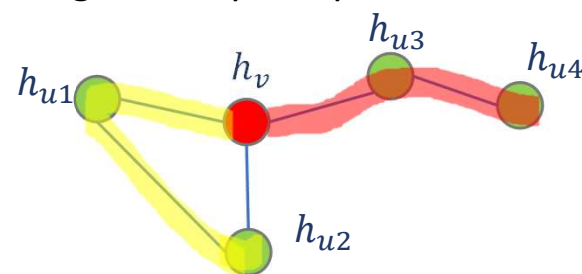
Stage I: Node Representation



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

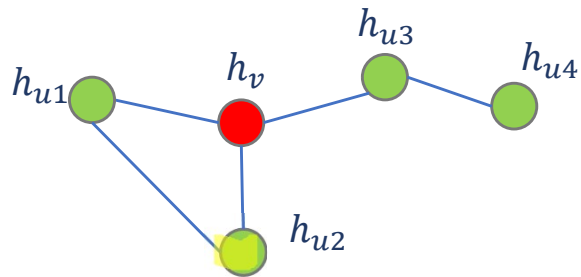
Equivalent to a GNN that needs no training



# PreGNN: Node- and Graph-level Pretraining

- (Hu et al. 2020)

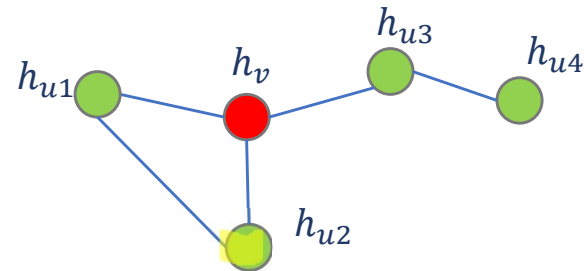
Stage I: Node Representation



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:

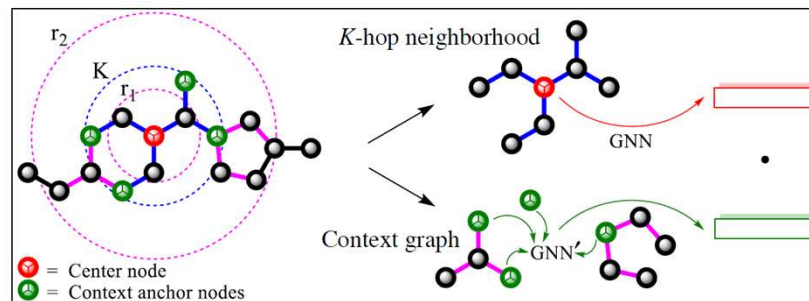
$$\min -\log\left(\sigma\left((h_v^K)^T C_v^G\right)\right) - I(v \neq v') \log(1 - \sigma\left((h_v^K)^T C_{v'}^{G'}\right))$$

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

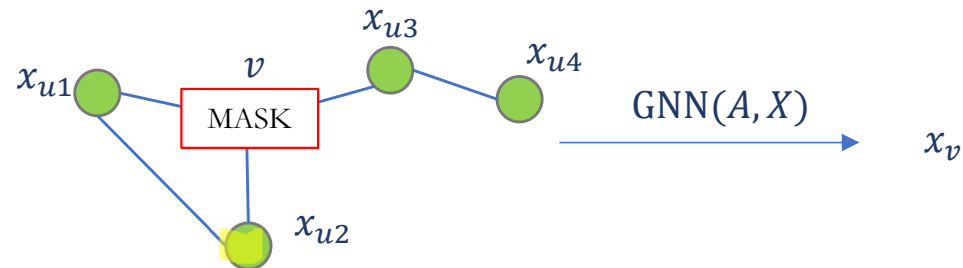
Hu et al. Strategies for Pre-training Graph Neural Networks.

# PreGNN

- (Hu et al. 2020)

Stage I: Node Representation

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



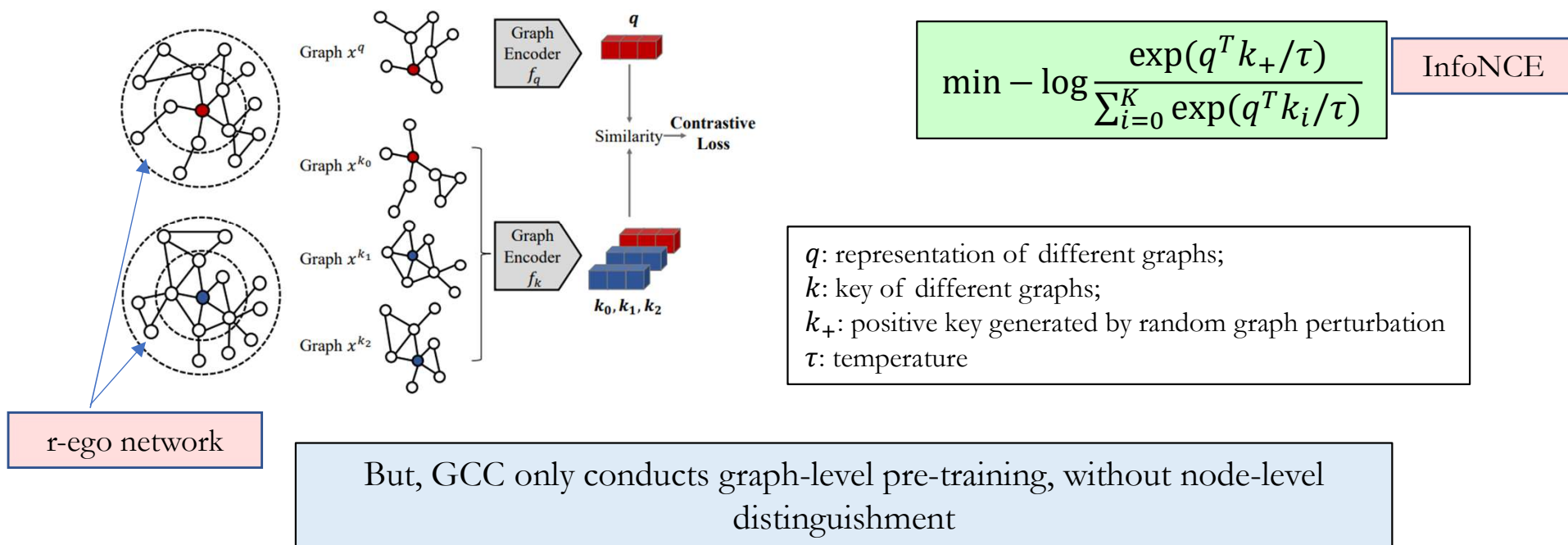
*Context Prediction  
Attribute Masking*

Hu et al. Strategies for Pre-training Graph Neural Networks.

# GCC: Contrastive learning

- (Qiu et al. 2020)

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



Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training.

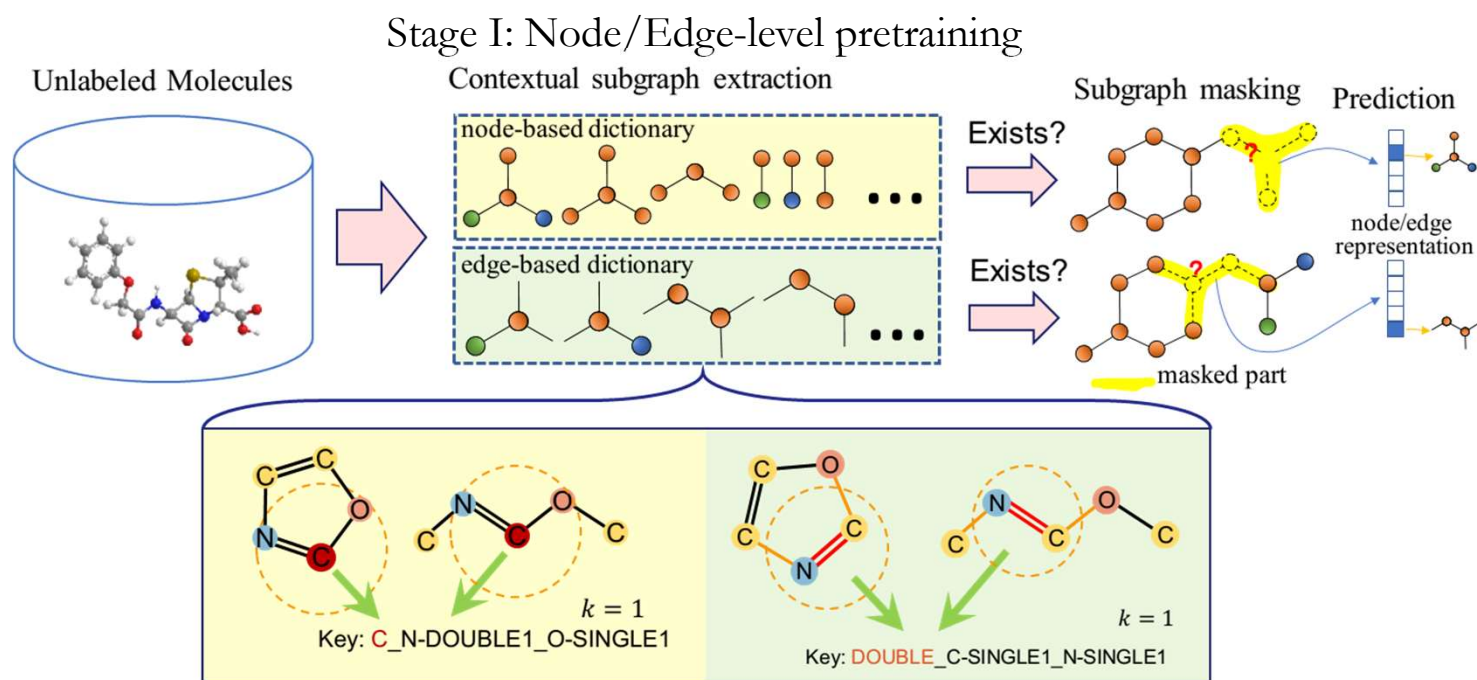
# GROVER (Rong et al. 2020)

---

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

# GROVER

- (Rong et al. 2020)



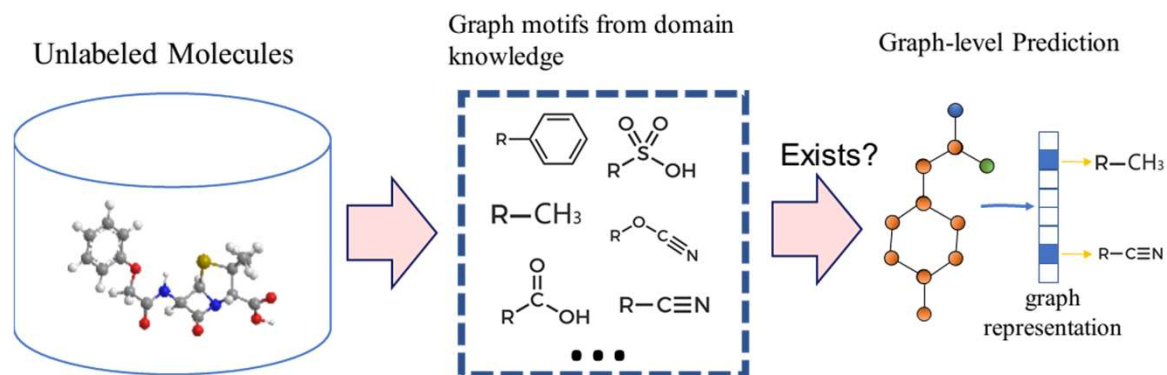
Predicting node/edge contexts instead of node labels can better capture local structures (multi-label)

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

# GROVER

- (Rong et al. 2020)

## Stage II: Graph-level pretraining



Predicting a graph if contains pre-defined graph motifs.

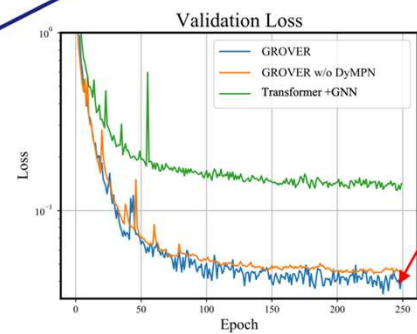
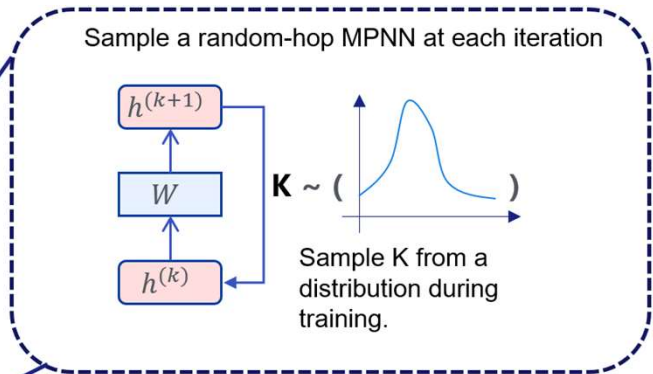
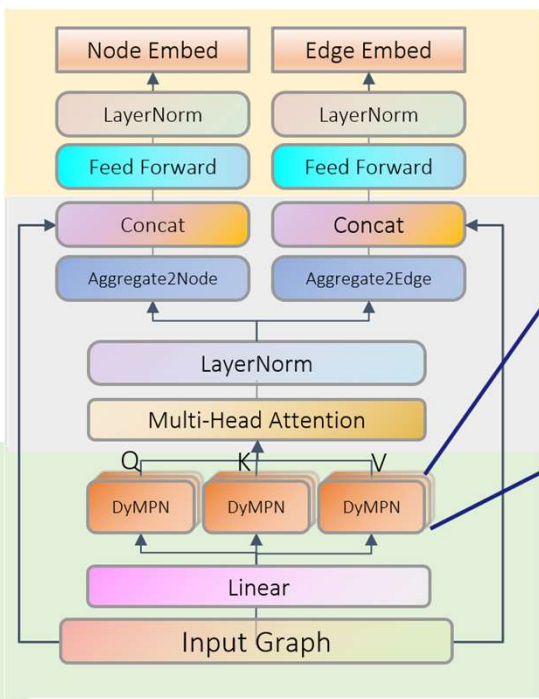
Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

# GROVER

- One more thing: GTransformer

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

- Output for both node embedding and edge embedding.
- Multi-Head Attention: model **global interaction** between nodes/edges.
- Long-range Residual Connection: alleviating the vanishing gradient and over-smoothing.
- MPNN: Extract **local structural information** of graphs.
- dyMPN: Randomize the message passing hops for the dynamic receptive field modeling.



Better generalization ability

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.



---

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

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

# GROVER

## Molecular classification

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

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

# 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]

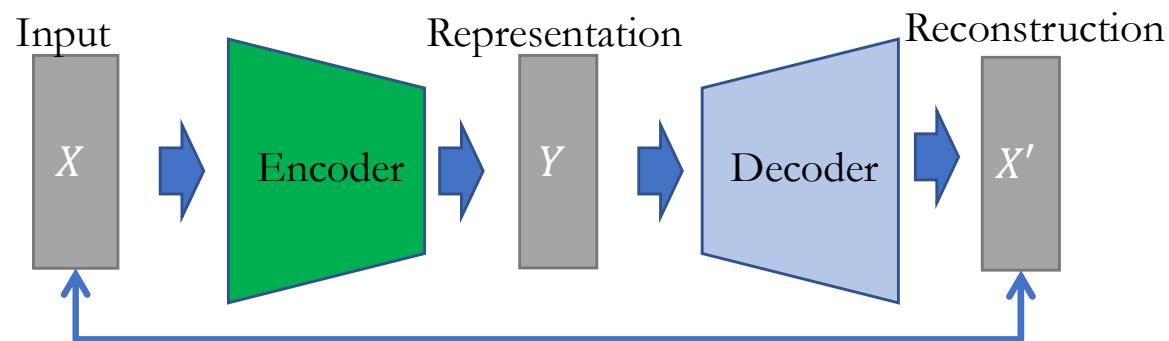


- [1] Kipf & Welling 2016; [2] Durán & Niepert 2017; [3] Hamilton et al. 2017;  
 [4] Liu et al. 2019; [5] Hu et al. 2020; [6] Qiu et al. 2020; [7] Rong et al. 2020;  
 [8] Veličković et al. 2019; [9] Peng et al. 2020; [10] Sun et al. 2020  
 [11] Hwang, Dasol, et al. 2020 [12] Peng, Zhen, et al.  
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

# What makes a good representation?

---

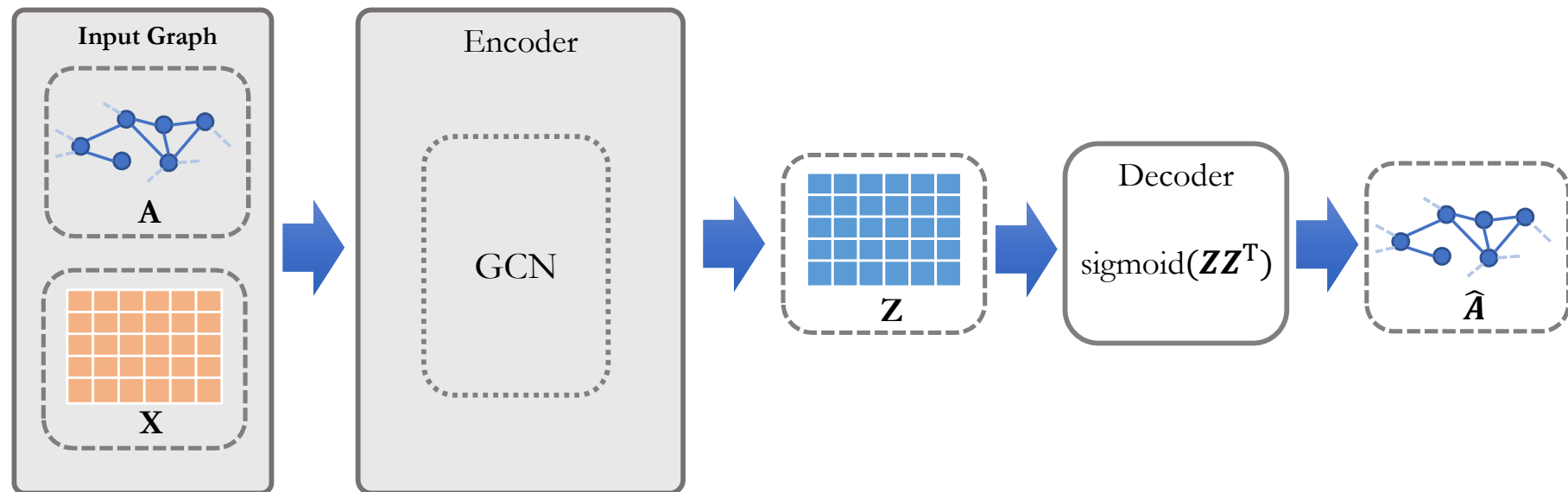
Auto-Encoder (AE)



“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)



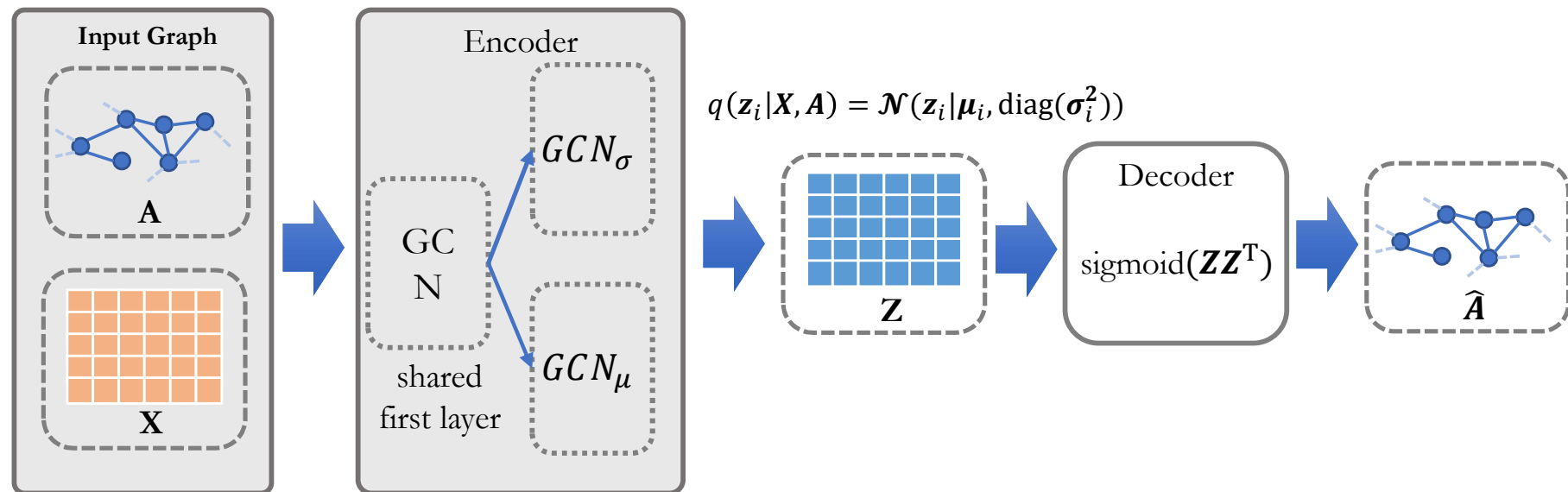
The overall loss:

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

The reconstruction loss.

Kipf, T. N., & Welling, M. (2016).

# Variational Graph Auto-Encoders (VGAE)



The overall loss:

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

The reconstruction loss.

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

Kipf, T. N., & Welling, M. (2016).

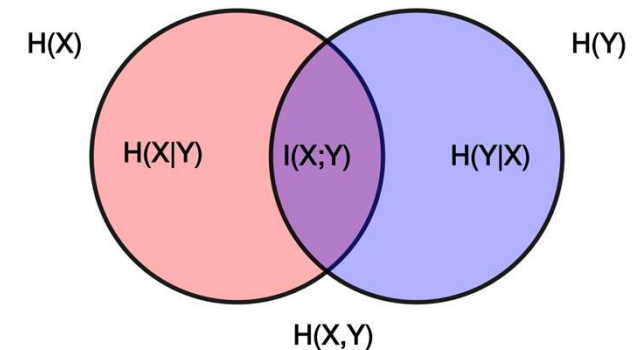
# What makes a good representation?

---

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

$$I(X; Y) = D_{KL}(p(X)p(Y) || p(X, Y)) \\ = H(X) - H(X|Y)$$

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

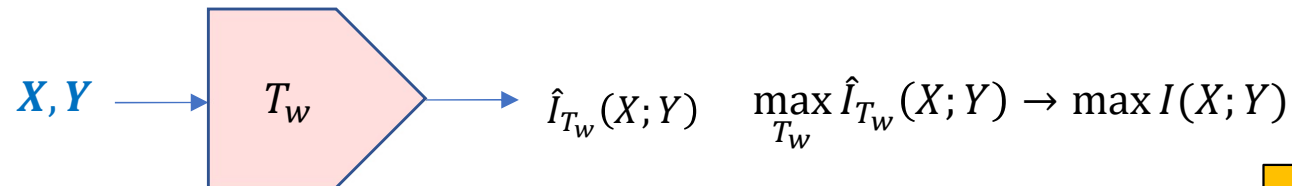
↑  
Reconstruction error

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):

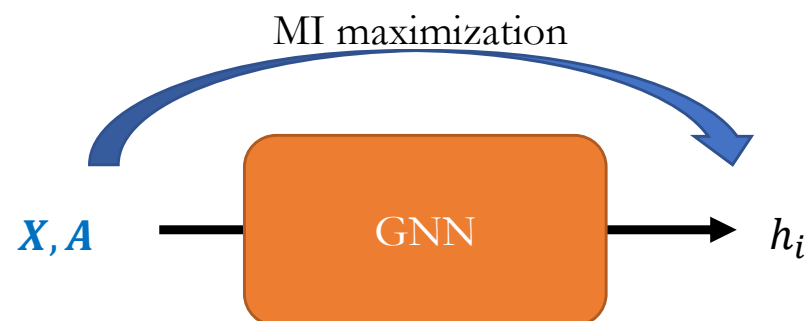
$$I^{\text{JSD}}(X; Y) \triangleq E_{p(X,Y)}[\log \sigma(T_w(x, y))] + E_{p(X)p(Y)}[\log(1 - \sigma(T_w(x, y)))]$$

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

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

# Deep Graph Infomax (DGI)

- (Velickovic et al. 2019)



The JSD MI estimator is applied:

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

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

$\tilde{h}_i$  negative sample

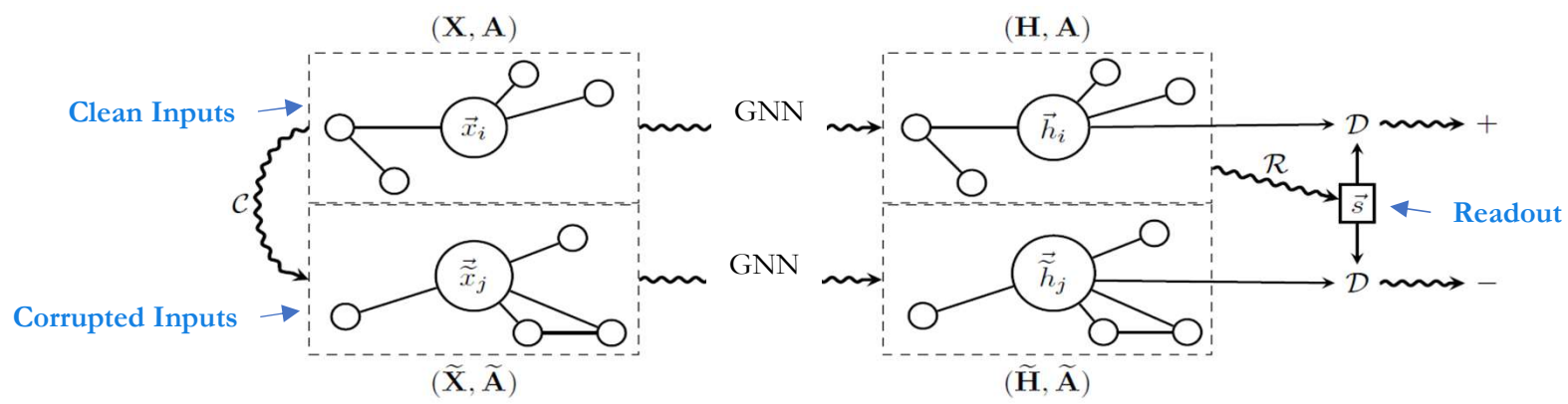
# Deep Graph Infomax (DGI)

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

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



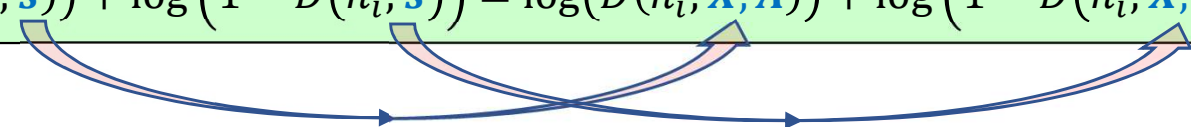
$$\max_{\text{GNN}} I(X, A; h_i) = \max_{\text{GNN}} \log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s}))$$



# 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,

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

---

- Some issues in DGI

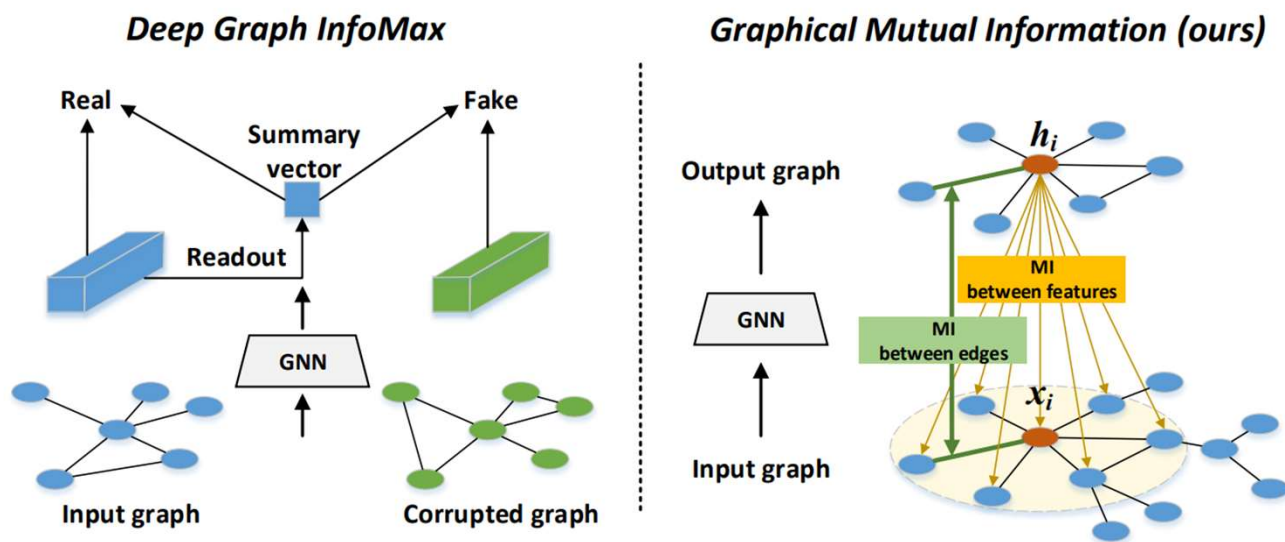
- Computing MI requires the injectivity of readout function
- It resorts to graph corruption to generate negative samples
- Distinct encoders and corruption functions for different tasks

*Indirect*

*Inefficient*

# GMI: Graphical Mutual Information

- (Peng et al. 2020)



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

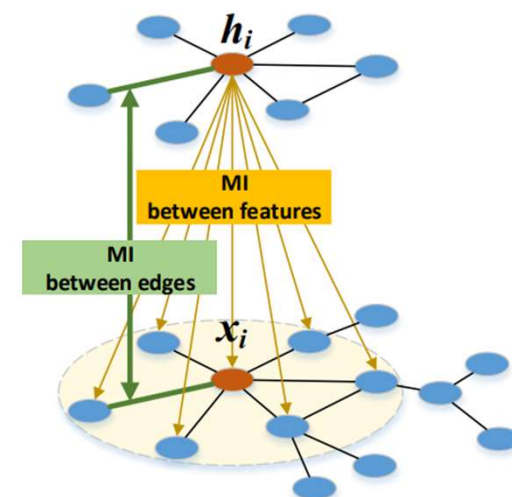
Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

# GMI: Graphical Mutual Information

We define that,

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

- 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)

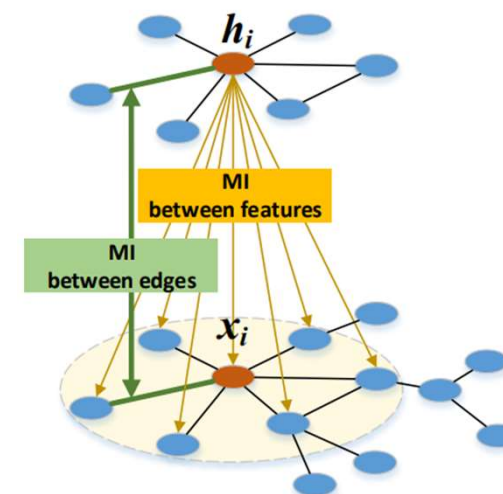
It can be proved that, if certain mild condition meets,

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

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 ± 0.1	69.3 ± 0.2	78.6 ± 0.2	93.8 ± 0.03	61.8 ± 0.04
DGI loss	82.2 ± 0.2	72.2 ± 0.2	78.9 ± 0.3	94.3 ± 0.02	62.3 ± 0.02
<b>FMI (ours)</b>	78.3 ± 0.1	72.0 ± 0.2	<b>79.1 ± 0.3</b>	<b>94.7 ± 0.03</b>	<b>64.8 ± 0.03</b>
<b>GMI-mean (ours)</b>	<b>82.7 ± 0.1</b>	<b>73.0 ± 0.3</b>	<b>80.1 ± 0.2</b>	<b>95.0 ± 0.02</b>	<b>65.0 ± 0.02</b>
<b>GMI-adaptive (ours)</b>	<b>83.0 ± 0.3</b>	<b>72.4 ± 0.1</b>	<b>79.9 ± 0.2</b>	<b>94.9 ± 0.02</b>	<b>64.6 ± 0.03</b>

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

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

# GMI: Graphical Mutual Information

---

- Link Prediction

We use an universal backbone (GCN) for all tasks

Algorithm	Cora			BlogCatalog			Flickr			PPI
	20.0%	50.0%	70.0%	20.0%	50.0%	70.0%	20.0%	50.0%	70.0%	22.7%
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
FMI (ours)	<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>	75.1±0.2	<b>92.7±0.3</b>	<b>92.2±0.3</b>	<b>90.6±0.4</b>	<b>79.8±0.2</b>
GMI (ours)	<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>

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

# Summary

---

	Node-Classification	Link/Metapath Prediction	Graph-Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	$S^2$ 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]

# Applications

---

- **GNN in Social Networks**

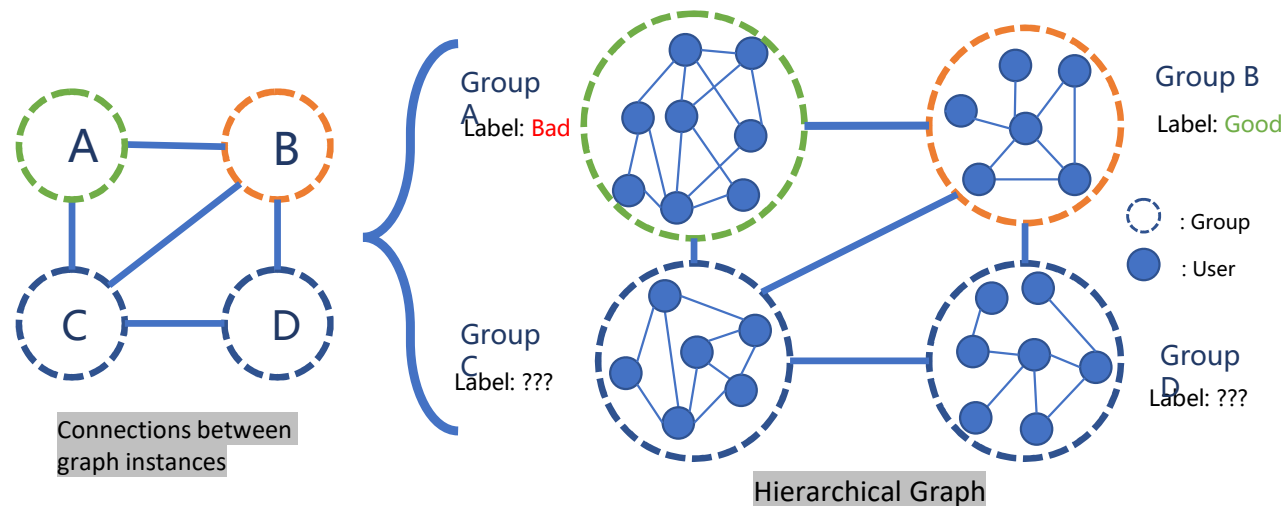
# GNN in Social Networks

---

- "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.

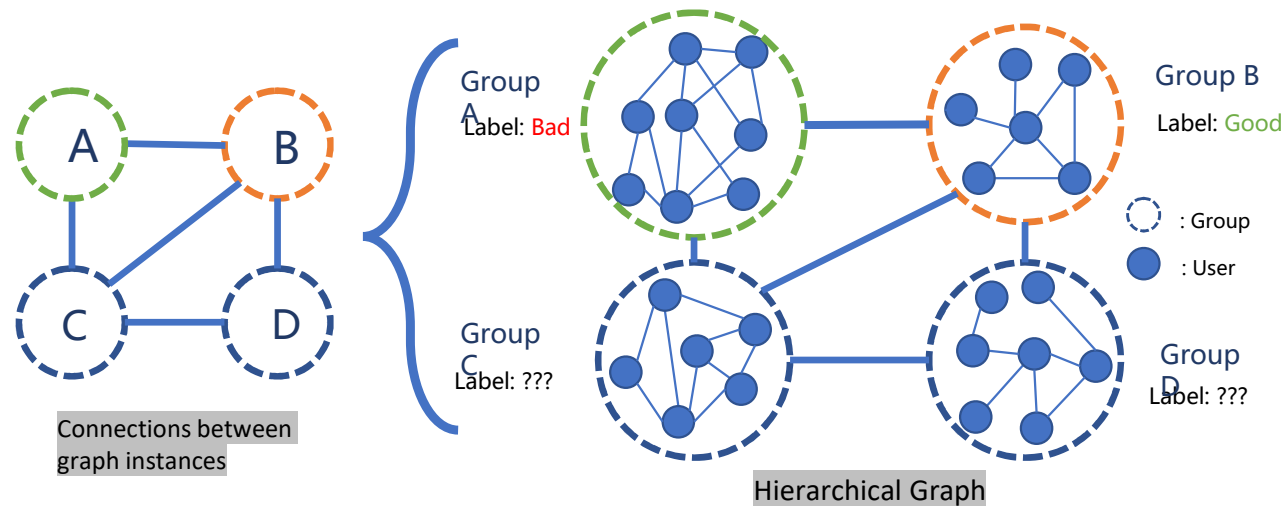


- **The Problem:** predicts the class label of graph instances in a hierarchical graph.

Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

# Hierarchical Graph Classification

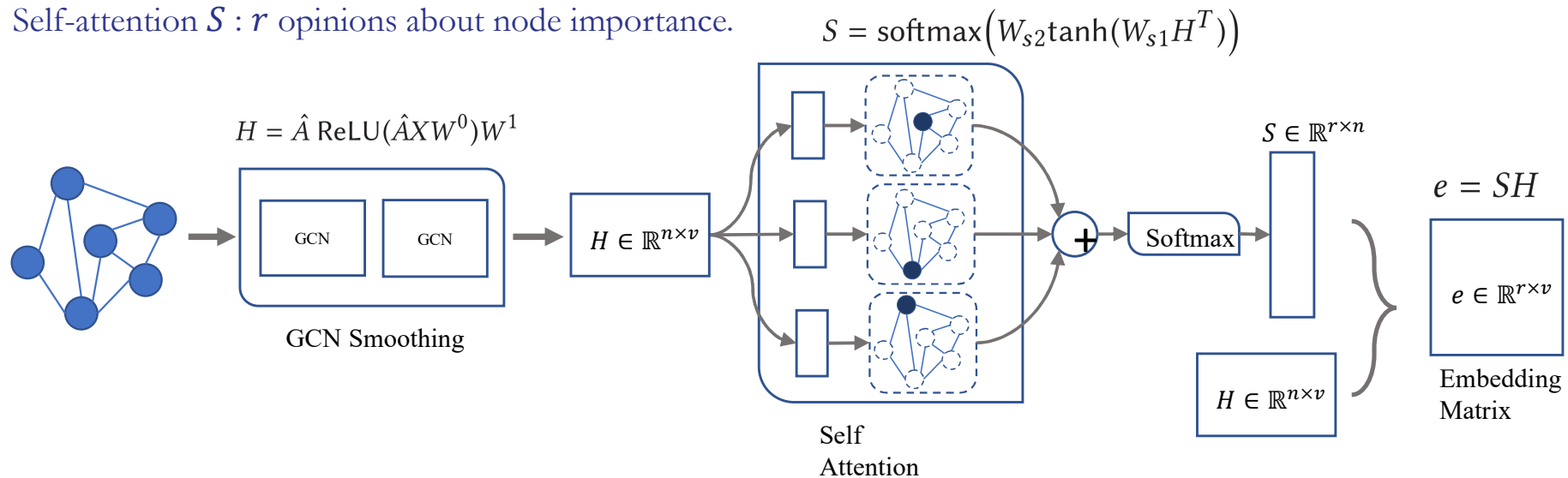
- **The Problem:** predicts the class label of graph instances in a hierarchical graph.
- **Challenges:**
  - How to represent the graphs with arbitrary size into a fixed-length vector?
  - How to incorporate the information of instance level and hierarchical level?



Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

# 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.

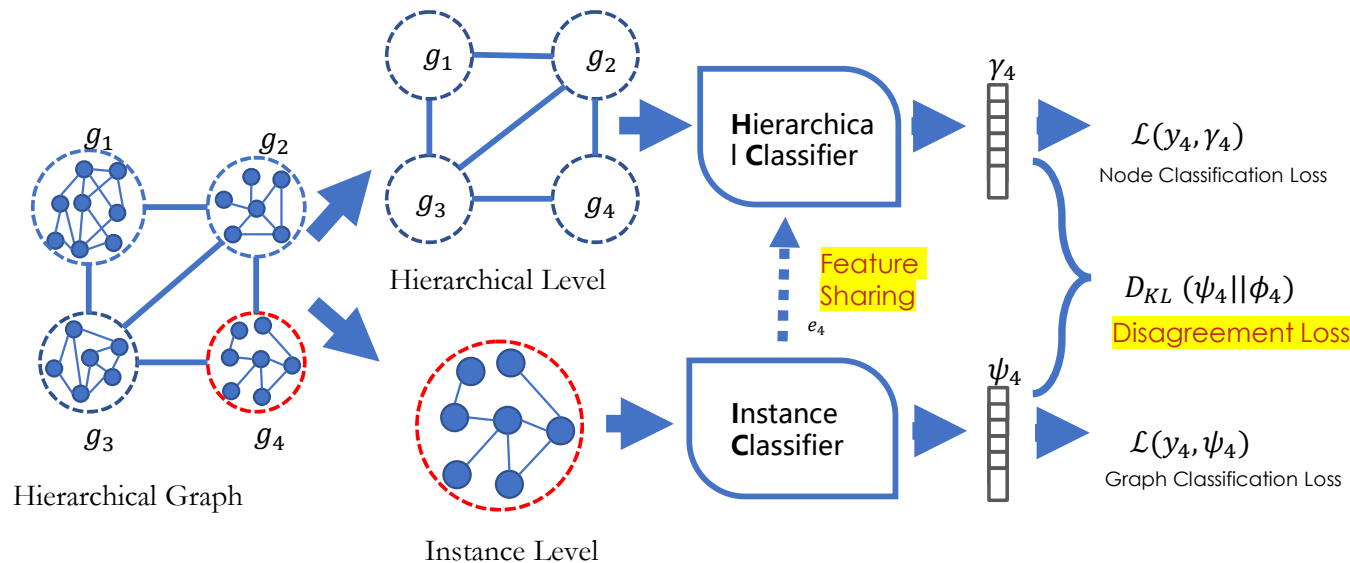


Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."



# 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:

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

Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

# 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

# GNN in Medical Imaging

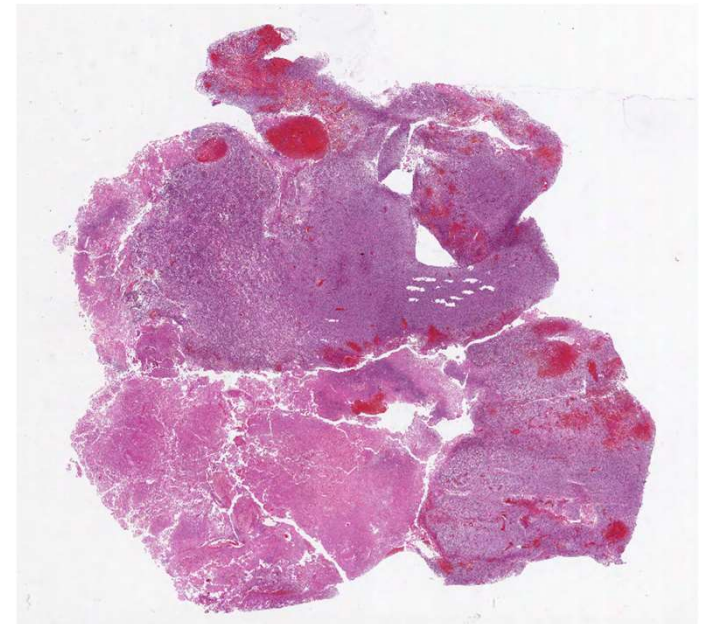
---

- **Survival Prediction**

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

- **Whole Slide Images**

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



# GNN in Medical Imaging

---

- Cox proportional hazard function

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

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

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

where,  $Y$  is the observation time.

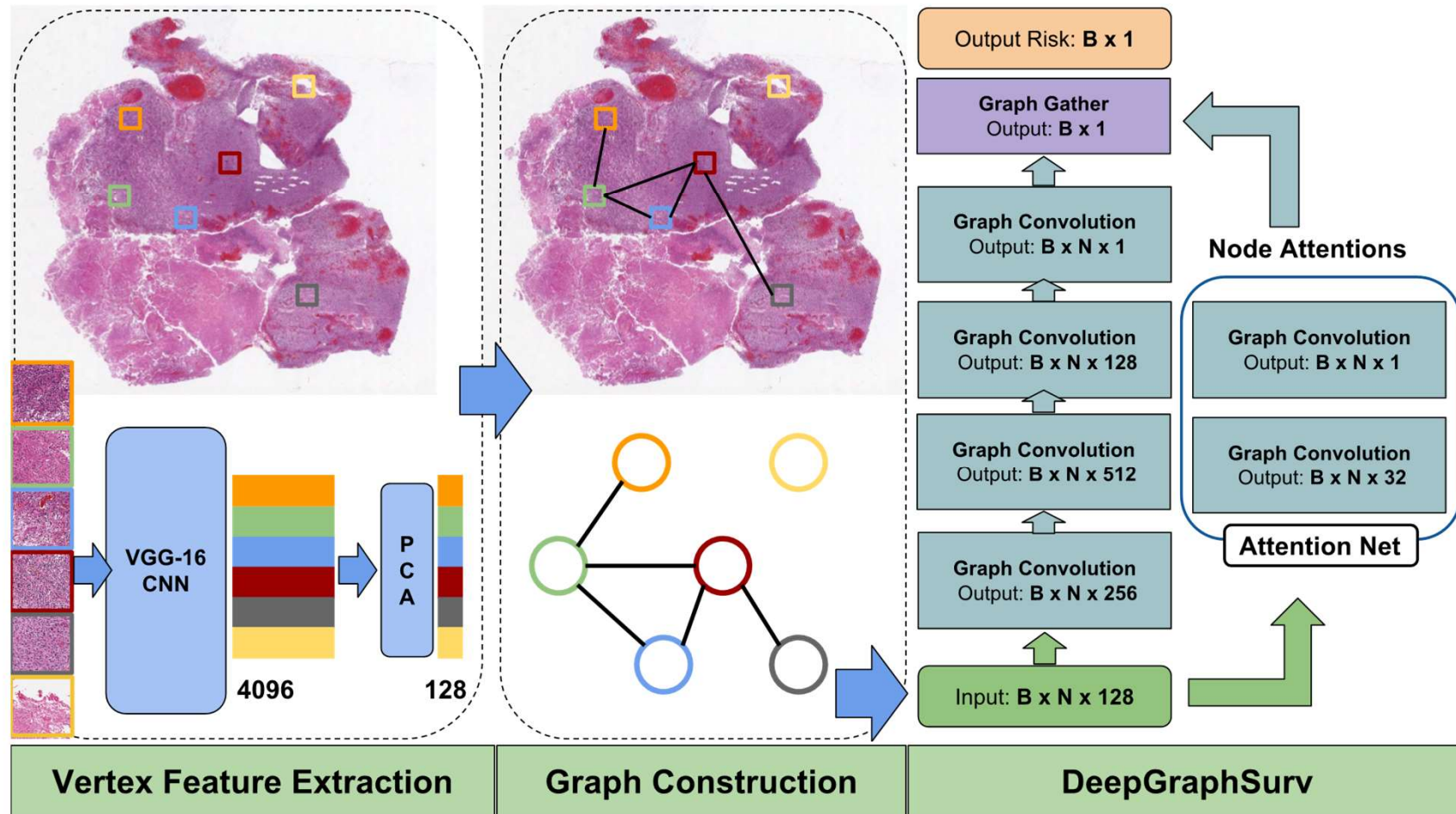
$$\theta_j = \exp(X_j \cdot \beta)$$

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

- Log likelihood as object function:

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

# GNN in Medical Imaging



# GNN in Medical Imaging

---

- **Pathological Images and Patient Survival Time and Label**
  - TCGA, The Cancer Genome Atlas
  - NLST, National Lung Screening Trials

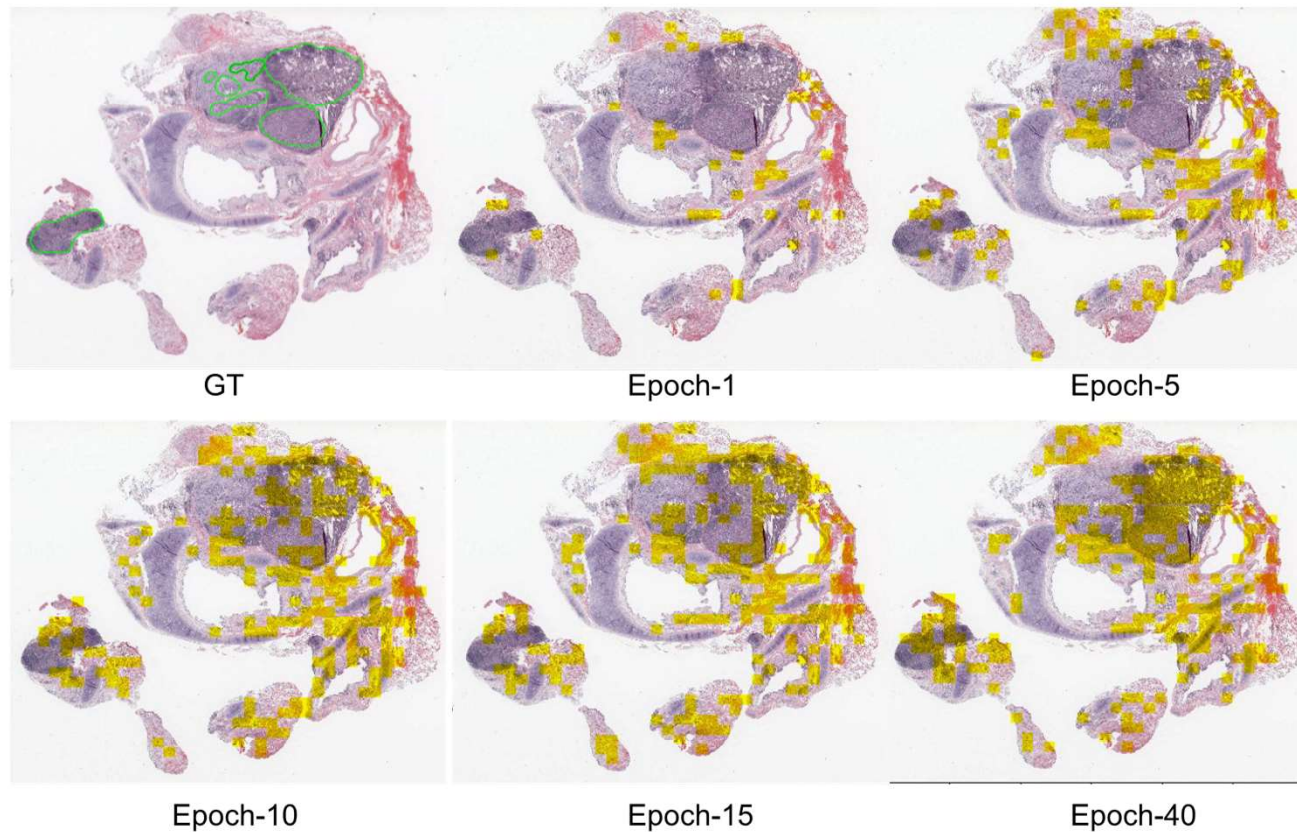
Database	Cancer Subtype	No. Patient	No. WSI	Quality	Avg. Size
TCGA	LUSC	463	535	medium	0.72 GB
TCGA	GBM	365	491	low	0.50 GB
NLST	ADC & SCC	263	425	high	0.74 GB

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

# GNN in Medical Imaging

---

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





# GNN in Medical Imaging

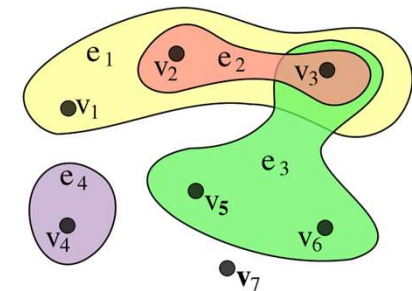
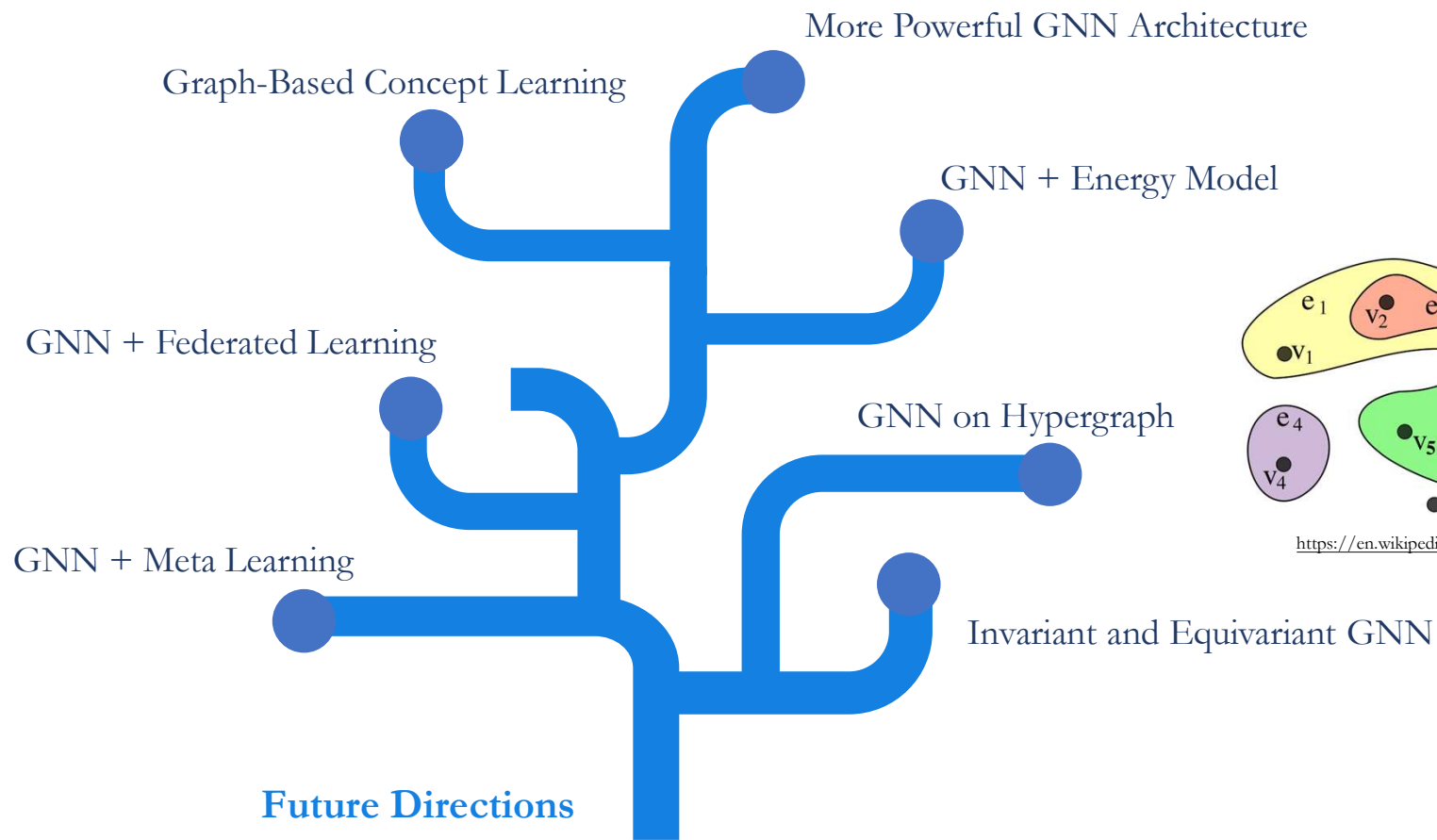
---

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>

---

# Bibliography

# Bibliography

---

- Sperduti, Alessandro, and Antonina Starita. "Supervised neural networks for the classification of structures." *IEEE Transactions on Neural Networks* 8.3 (1997): 714-735.
- 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.
- Scarselli, Franco, et al. "The graph neural network model." *IEEE Transactions on Neural Networks* 20.1 (2008): 61-80.
- Li, Yujia, et al. "Gated graph sequence neural networks." arXiv preprint arXiv:1511.05493 (2015).
- Defferrard, Michaël, Xavier Bresson, and Pierre Vandergheynst. "Convolutional neural networks on graphs with fast localized spectral filtering." *Advances in neural information processing systems*. 2016.
- 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).
- Bruna, Joan, et al. "Spectral networks and locally connected networks on graphs." arXiv preprint arXiv:1312.6203 (2013).
- Niepert, Mathias, Mohamed Ahmed, and Konstantin Kutzkov. "Learning convolutional neural networks for graphs." *International conference on machine learning*. 2016.
- Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907 (2016).
- Xu, Bingbing, et al. "Graph Wavelet Neural Network." *International Conference on Learning Representations*. 2018.
- Chami, Ines, et al. "Hyperbolic graph convolutional neural networks." *Advances in neural information processing systems*. 2019.
- Liao, Renjie, et al. "LanczosNet: Multi-Scale Deep Graph Convolutional Networks." *International Conference on Learning Representations*. 2018.
- Veličković, Petar, et al. "Graph Attention Networks." *International Conference on Learning Representations*. 2018.
- Zhang, Jiani, et al. "Gaan: Gated attention networks for learning on large and spatiotemporal graphs." arXiv preprint arXiv:1803.07294 (2018).
- Chang, Heng, et al. "Spectral Graph Attention Network." arXiv preprint arXiv:2003.07450 (2020).
- Ying, Zhitaο, et al. "Hierarchical graph representation learning with differentiable pooling." *Advances in neural information processing systems*. 2018.
- Ma, Yao, et al. "Graph convolutional networks with eigenpooling." *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2019.
- Atwood, James, and Don Towsley. "Diffusion-convolutional neural networks." *Advances in neural information processing systems*. 2016.
- Abu-El-Haija, Sami, et al. "MixHop: Higher-Order Graph Convolutional Architectures via Sparsified Neighborhood Mixing." *International Conference on Machine Learning*. 2019.
- Klicpera, Johannes, Aleksandar Bojchevski, and Stephan Günnemann. "Predict then Propagate: Graph Neural Networks meet Personalized PageRank." *International Conference on Learning Representations*. 2018.
- Gilmer, Justin, et al. "Neural Message Passing for Quantum Chemistry." *ICML*. 2017.
- Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective." *The World Wide Web Conference*. 2019.

# Bibliography

---

- Hamilton, Will, Zhitaoying, and Jure Leskovec. "Inductive representation learning on large graphs." Advances in neural information processing systems. 2017.
- 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.
- Huang, Wenbing, et al. "Adaptive sampling towards fast graph representation learning." Advances in neural information processing systems. 2018.
- 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.
- Zeng, Hanqing, et al. "GraphSAINT: Graph Sampling Based Inductive Learning Method." International Conference on Learning Representations. 2020.
- Morris et al. Weisfeiler and Leman Go Neural Higher-order Graph Neural Networks. AAAI, 2019.
- Xu et al. How Powerful Are Graph Neural Networks. ICLR, 2019.
- Maron et al. Invariant and Equivariant Graph Networks. ICLR, 2019.
- Maron et al. On the Universality of Invariant Networks. ICML, 2019.
- Maron et al. Provably Powerful Graph Networks, NeurIPS. 2019.
- Dehmamy et al. Understanding the Representation Power of Graph Neural Networks in Learning Graph Topology. NeurIPS, 2019.
- Sato et al. Approximation Ratios of Graph Neural Networks for Combinatorial Problems. NeurIPS, 2019.
- Loukas, What graph neural networks cannot learn: depth vs width. ICLR, 2020.
- Garg et al. Generalization and Representational Limits of Graph Neural Networks. ICML, 2020.
- Shervashidze et al. Weisfeiler-Lehman Graph Kernels. JRML, 2011.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems, 2(4):303–314, 1989.
- Hornik, K. Approximation capabilities of multilayer feedforward networks. Neural networks, 4(2):251–257, 1991.
- Chen et al. On the equivalence between graph isomorphism testing and function approximation with GNNs. NeurIPS, 2019.
- Kipf & Welling. Variational Graph Auto-Encoders. arXiv, 2016.
- Durán & Niepert. Learning Graph Representations with Embedding Propagation. NeurIPS, 2017.
- Liu et al. N-Gram Graph: Simple Unsupervised Representation for Graphs, with Applications to Molecules. NeurIPS, 2019.
- Hu et al. Strategies for Pre-training Graph Neural Networks. ICLR, 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.
- Veličković et al. Deep Graph Infomax. ICLR, 2019.
- Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization. WWW, 2020.
- Sun et al. InfoGraph: Unsupervised and Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization. ICLR, 2020.

# Bibliography

---

- Hinton, G. E., & Salakhutdinov, R. R. Reducing the Dimensionality of Data with Neural Networks. *Science*, 2006.
- Vincent et al. Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion. *JMLR*, 2010.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *NeurIPS*, 2018.
- Belghazi et al. Mine: mutual information neural estimation. *ICML*, 2018.
- Nowozin et al. f-gan: Training generative neural samplers using variational divergence minimization. *NeurIPS*, 2016.
- Hjelm et al. Learning deep representations by mutual information estimation and maximization. *ICLR*, 2019.
- Wang, X., & Gupta, A. Videos as Space-Time Region Graphs. *ECCV*, 2018.
- Yan et al. Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. *AAAI*, 2018.
- Zeng, et al. Graph convolutional networks for temporal action localization. *ICCV*, 2019
- Bian, Tian, et al. "Rumor Detection on Social Media with Bi-Directional Graph Convolutional Networks." *AAAI* 2020.
- Junchi, Yu, et al. "Graph Information Bottleneck for Subgraph Recognition.", *ICLR* 2021.