

# Graph Neural Networks and Self-supervised Learning

Jie Tang

Computer Science  
Tsinghua University



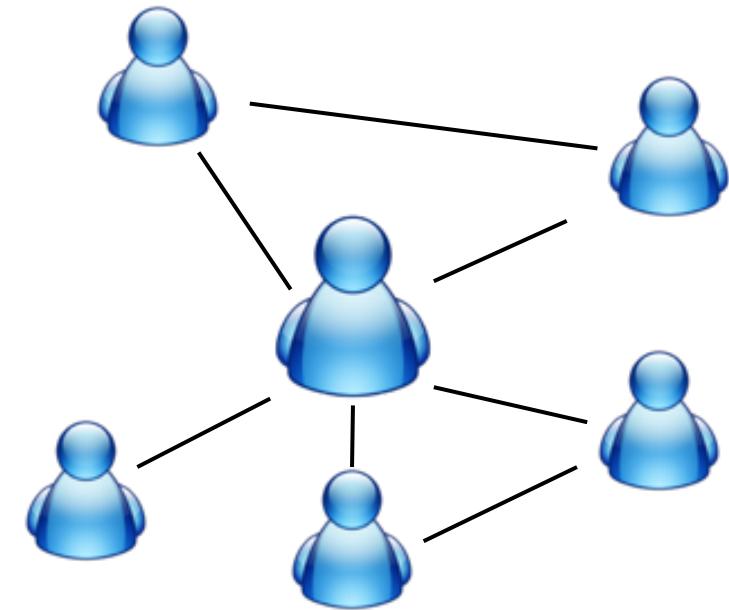
# A bit about Jie

- Jie Tang, IEEE Fellow, Professor, Associate Chair of Dept. of CS at Tsinghua University. Interests include social network, data mining, machine learning, knowledge graph.
  - SIGKDD Test of Time Award
  - SIGKDD Service Award
  - NSFC for Distinguished Young Scholars, CCF Young Scientist Award
  - 2<sup>nd</sup> National S&T Award
- PC Co-Chair
- Editor in Chief
- Have published over 200 papers in top journals and conferences, including KDD (30), IJCAI/AAAI (30), NIPS/ICML (10+), IEEE Trans. (30), Machine Learning J
- #Citations: 20,000+ and *h*-index: 74
- ArnetMiner <https://www.aminer.org/> for academic researcher network analysis, which attracted 20 million users from 220 countries/regions.
- HP: <http://keg.cs.tsinghua.edu.cn/jietang/>

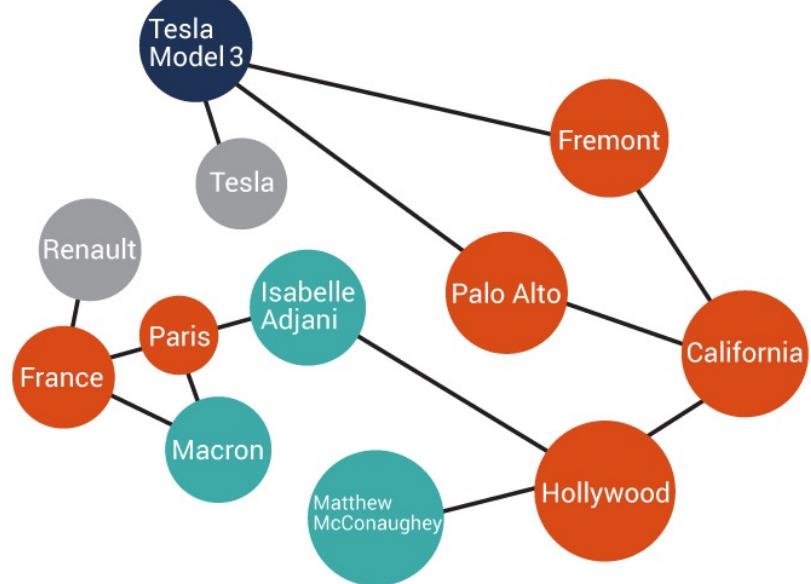


Work Hard and Play Harder!

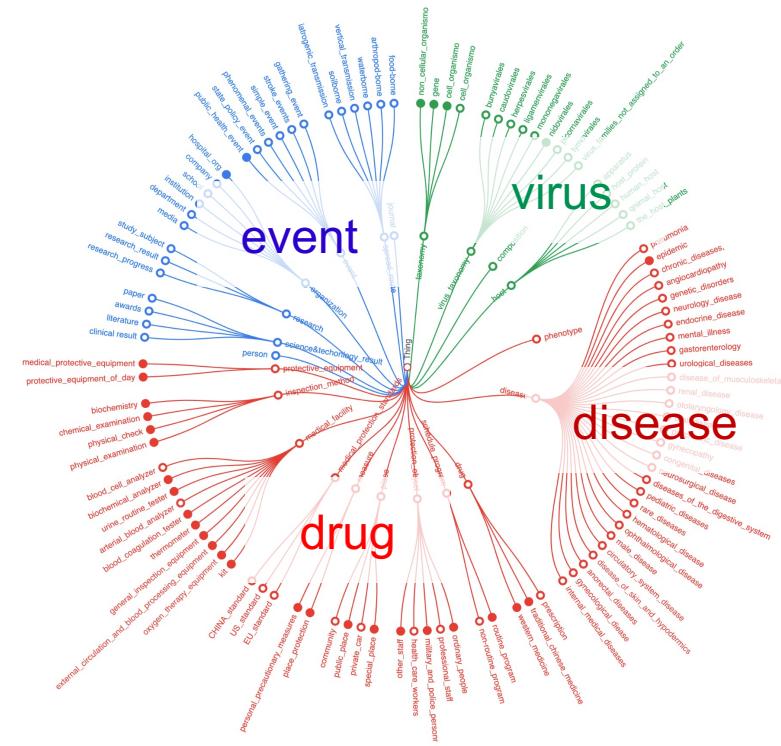
# Networked Data



# Social Network



# Knowledge Graph



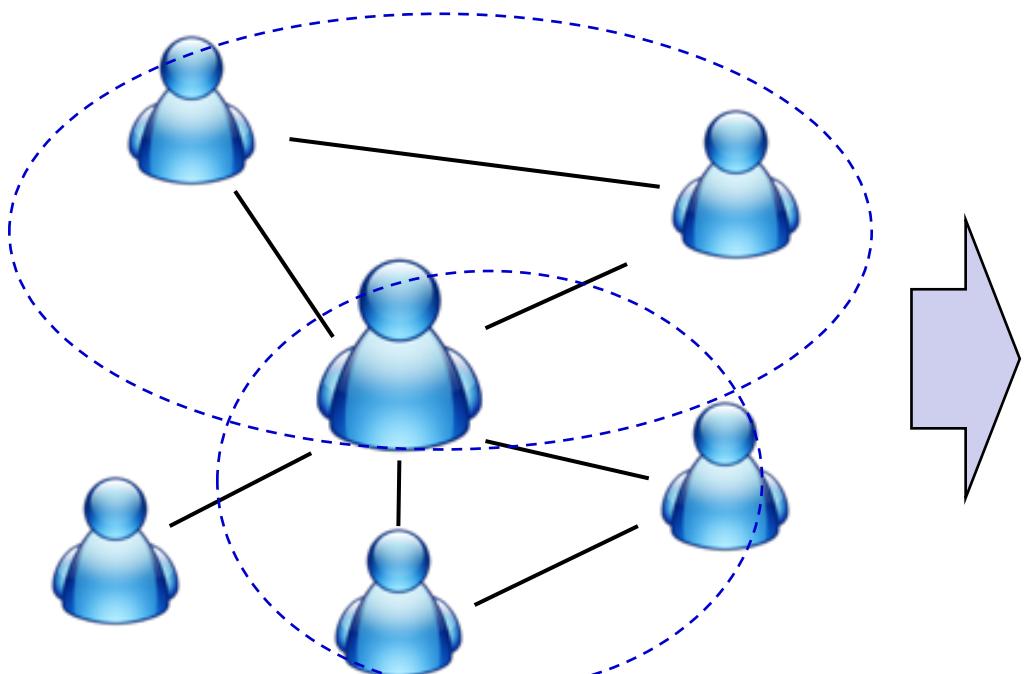
# COVID Graph

# Machine Learning with Networks

- ML tasks in networks:
  - Node classification
    - Predict a type of a given node
  - Link prediction
    - Predict whether two nodes are linked
  - Community detection
    - Identify densely linked clusters of nodes
  - Network similarity
    - How similar are two (sub)networks?

# Representation Learning on Networks

## Representation Learning/ Graph Embedding



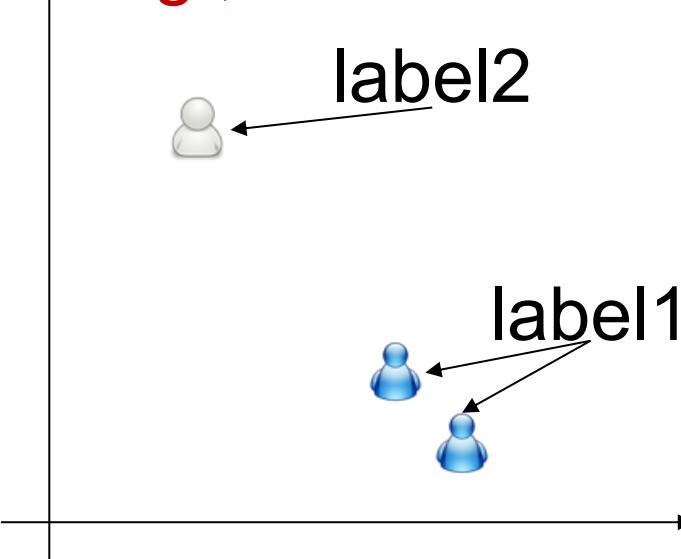
$d$ -dimensional vector,  $d \ll |V|$



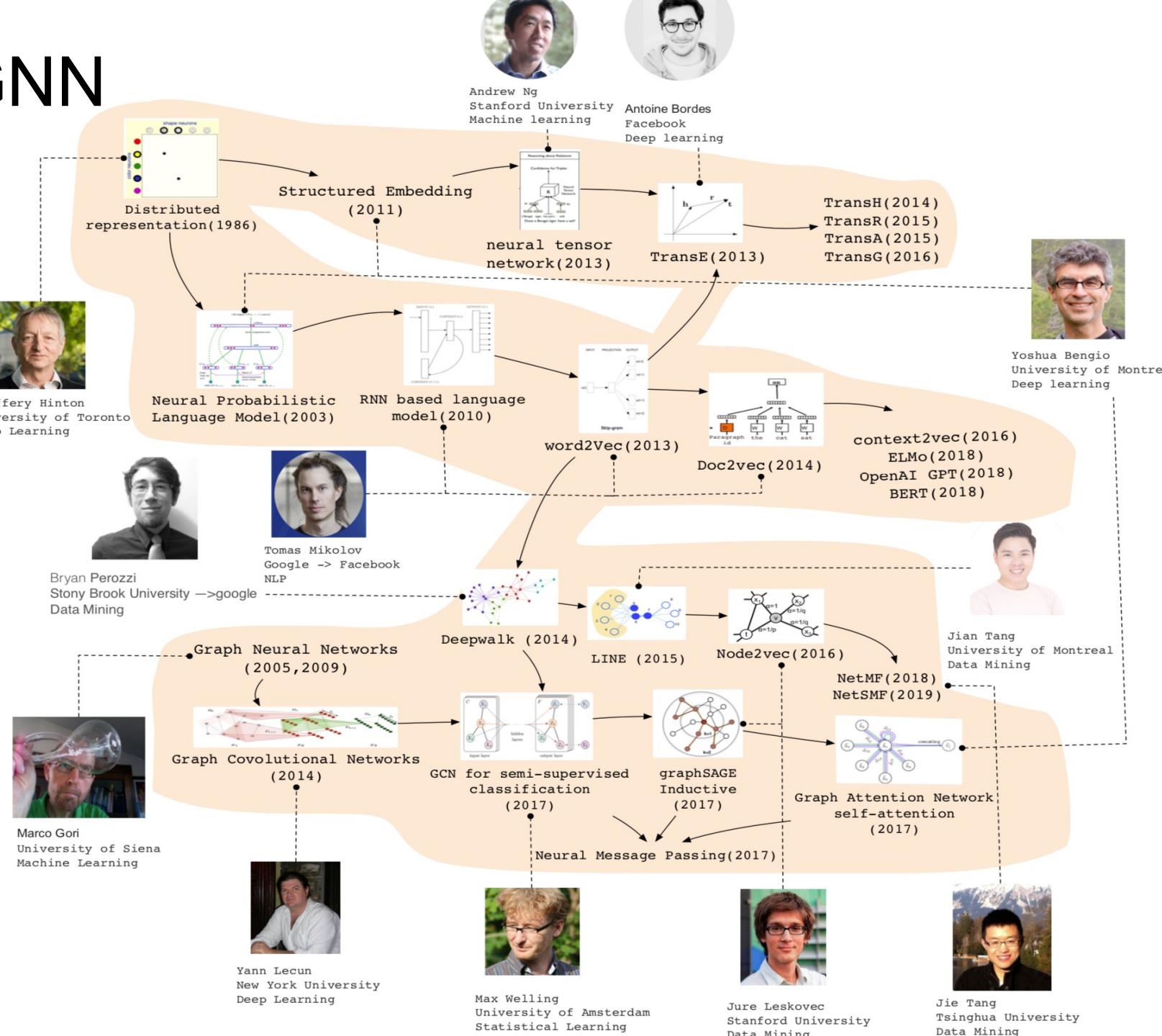
0.8	0.2	0.3	...	0.0	0.0
-----	-----	-----	-----	-----	-----

Users with the **same label** are located in **closer**

e.g., node classification



# GRL: NE&GNN



# 超大规模图神经网络

## ① 基于图神经网络的大规模图数据挖掘

现 实场景中充斥着大量图结构数据，涵盖社交网络、推荐系统、学术网络等，传统深度学习难以处理大规模图数据中复杂的结构化信息。清华大学计算机系知识工程实验室团队深入研究图神经网络，在图表示学习、图预训练模型、图神经网络鲁棒性、异构图学习等领域取得突破，为大规模图数据挖掘提供了理论依据和技术支持。

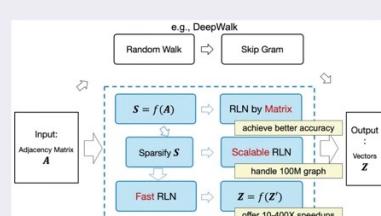
## ② 大规模图表示学习

如 何从大规模图数据中抽取出有效的表示是图数据挖掘的一大难题。以往方法如 DeepWalk、LINE、Node2Vec 计算复杂度高并缺乏理论支持，团队从理论上证明此类方法实际上可以统一成矩阵分解的近似形式，并提出基于显式矩阵分解的 NetMF 方法，效果有显著提升。NetMF 进一步将该方法扩展到稀疏图数据上，并在百万级节点规模的图数据实验中，将原本所需数月的计算缩短到一天以内。

团队进一步优化基于稀疏矩阵分解的图表示学习方法，提出引入频域图嵌入的 ProNE，在上亿级节点图上计算效率提高上百倍。

团队还将图表示学习拓展到异构图中，提出结合多模态异构图属性的 GATNE，在四亿商品网络上以极高精度预测节点间关系，性能较同质图方法显著提升。

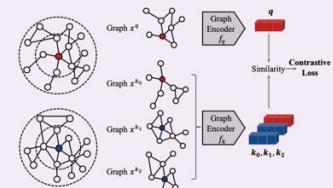
此外，为了更好地评价图神经网络在异构图上的性能，团队提出了异构图学习基准 HGB，大大推动了异构图神经网络的发展。



Algorithm	Closed Matrix Form
DeepWalk	$\log \left( \text{vol}(G) \left( \frac{1}{t} \sum_{i=1}^t (\mathbf{D}^{-1} \mathbf{A})^i \right) \mathbf{D}^{-1} \right) - \log b$
LINE	$\log \left( \text{vol}(G) \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right) - \log b$
PTE	$\log \left( \alpha \text{vol}(G_{\text{ww}}) (\mathbf{D}_{\text{row}}^{\text{ww}})^{-1} \mathbf{A}_{\text{ww}} (\mathbf{D}_{\text{col}}^{\text{ww}})^{-1} + \beta \text{vol}(G_{\text{dw}}) (\mathbf{D}_{\text{row}}^{\text{dw}})^{-1} \mathbf{A}_{\text{dw}} (\mathbf{D}_{\text{col}}^{\text{dw}})^{-1} + \gamma \text{vol}(G_{\text{wd}}) (\mathbf{D}_{\text{row}}^{\text{wd}})^{-1} \mathbf{A}_{\text{wd}} (\mathbf{D}_{\text{col}}^{\text{wd}})^{-1} \right) - \log b$
node2vec	$\log \left( \frac{1}{2} \sum_{i=1}^T \left( \sum_u X_{u, i} P'_{i, u} + \sum_u X_{i, u} P'_{u, i} \right) \right) - \log b$

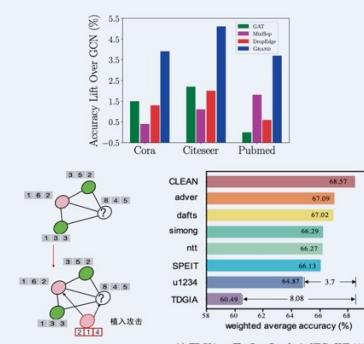
## ③ 大规模图预训练模型

图 数据由于其独特的结构化属性，难以进行迁移。团队提出图预训练技术 GCC，将图结构信息通过随机采样加对比学习的方式进行编码，训练出适用于不同类型、不同规模图数据的预训练模型。该模型在大规模图数据上有效提高节点分类、子图分类等下游任务上的表现。



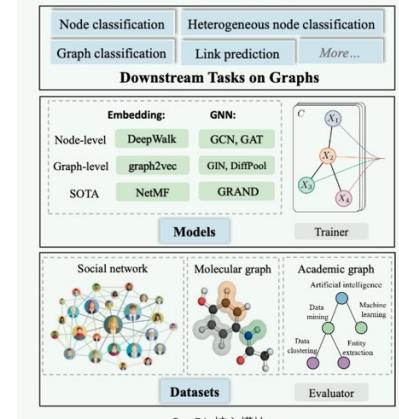
## ④ 图神经网络鲁棒性

图 神经网络容易受到对抗攻击风险，团队从攻防两方面展开鲁棒性研究。提出了随机图神经网络 GRAND，通过随机信息传播过程减轻攻击带来的影响；提出了拓扑图植入攻击 TDGIA，针对图数据中的薄弱连接，注入少量节点就能大幅降低图神经网络性能。

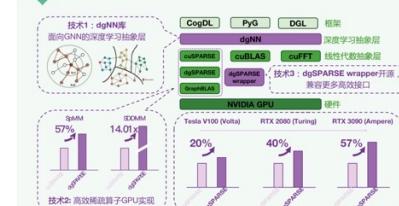


## ⑤ 高效易用的图深度学习框架 CogDL

为 了更好推进图神经网络发展，团队提出了 CogDL 图深度学习框架。CogDL 的核心理念包括易用接口、可复现性和高效计算，帮助研究人员可以轻松地训练和比较基线算法或自定义模型，以进行节点分类，链接预测，图分类，社区发现等任务。它提供了多种图学习算法实现，包括图神经网络、图表示学习、图预训练模型、图鲁棒性方法等，并基于高效稀疏算子将模型训练速度提高多个量级。自 2020 年 4 月正式启动以来，CogDL 已经发布了若干个版本，持续不断地进行改进和优化，助力于图深度学习的研究和应用。



## CogDL 图深度学习框架



# GRL: NE&GNN

Network  
Embedding

Matrix  
Factorization

Graph Neural  
Networks

GNN  
Pre-Training



<https://alchemy.tencent.com/>

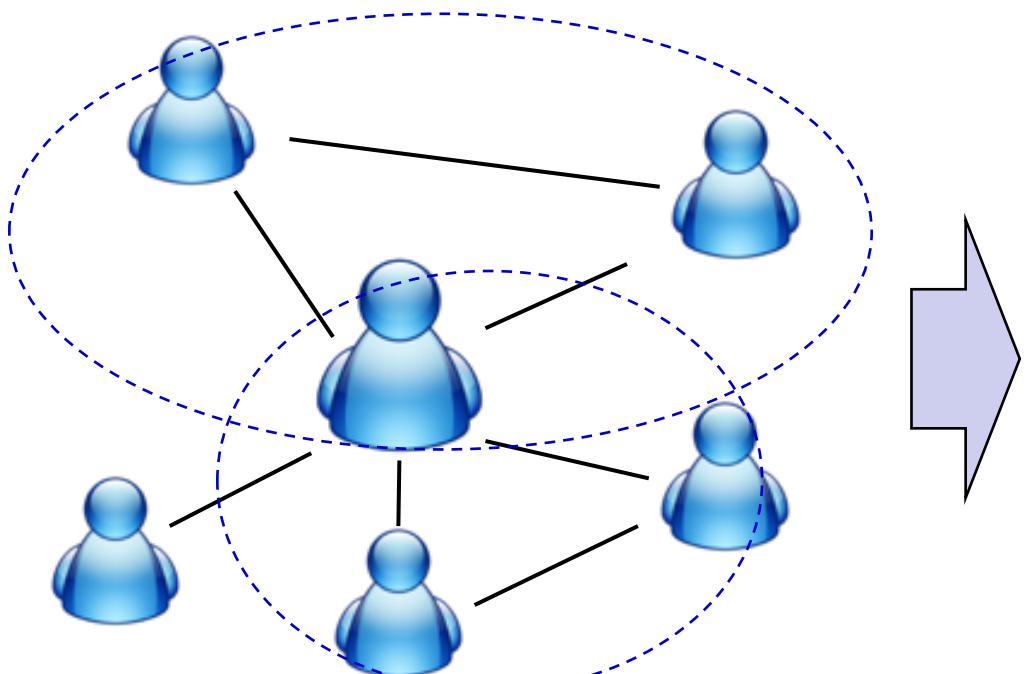


CogDL

<https://github.com/thudm/cogdl>

# Representation Learning on Networks

## Representation Learning/ Graph Embedding



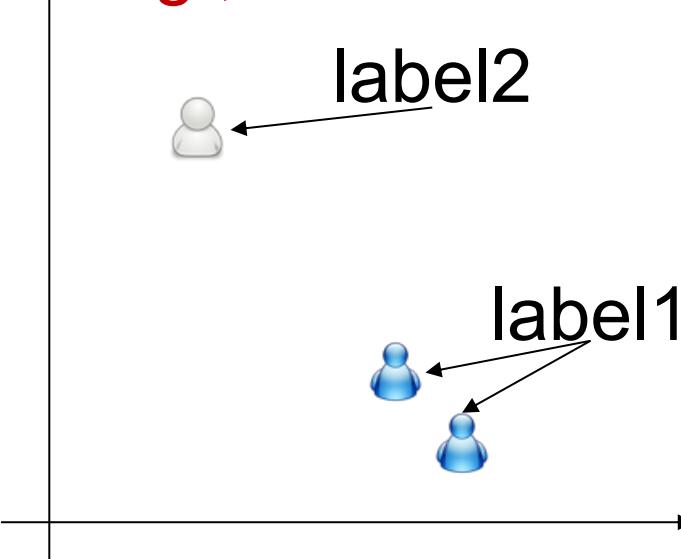
$d$ -dimensional vector,  $d \ll |V|$



0.8	0.2	0.3	...	0.0	0.0
-----	-----	-----	-----	-----	-----

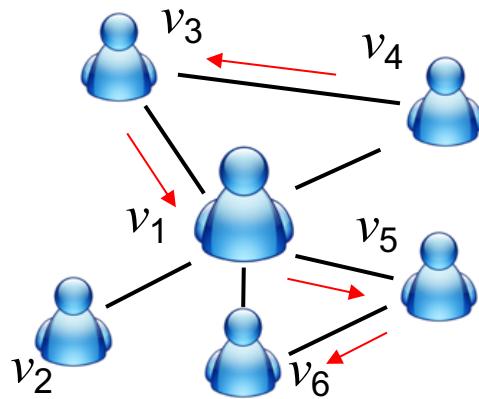
Users with the **same label** are located in **closer**

e.g., node classification

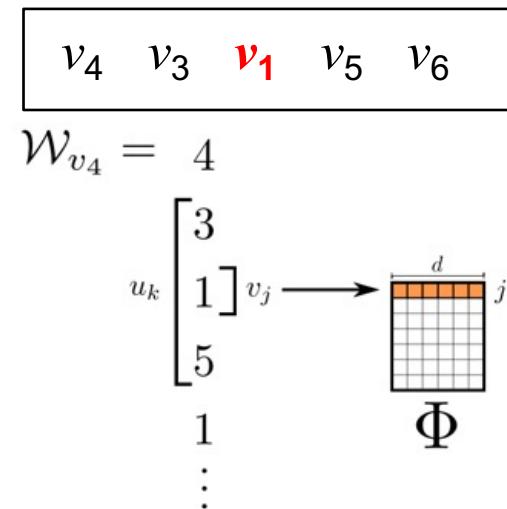


# DeepWalk: Random Walk + Word2Vec

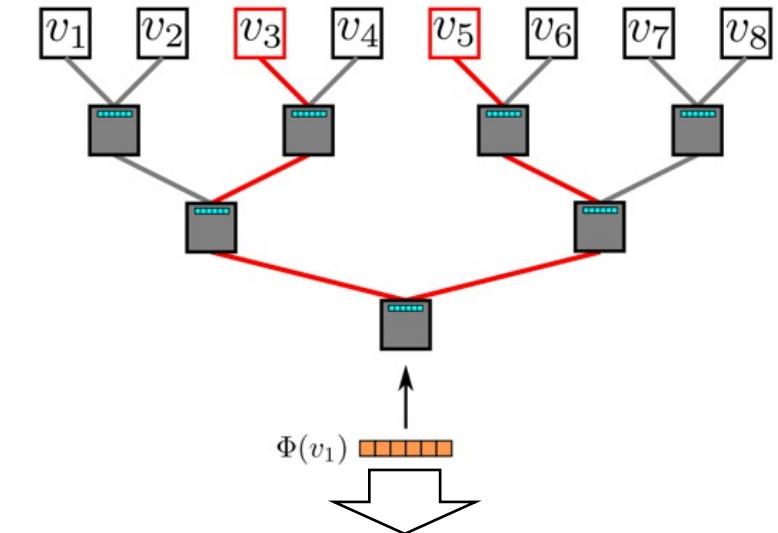
Random walk



One example RW path



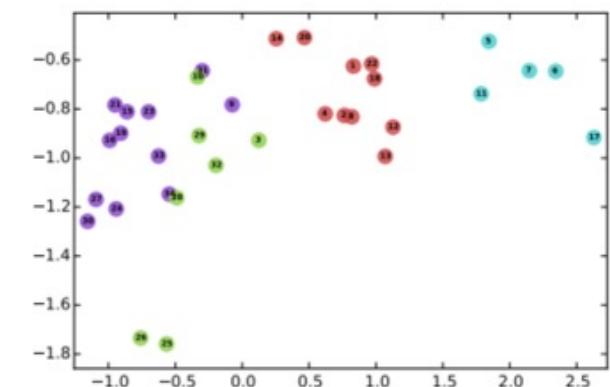
**SkipGram** with  
Hierarchical softmax



Hierarchical softmax

$$P(\{v_{i-w}, \dots, v_{i+w}\} \setminus v_i | \Phi(v_i)) = \prod_{j=i-w, j \neq i}^{i+w} P(v_j | \Phi(v_i))$$

$$P(v_j | \Phi(v_i)) = \prod_{l=1}^{\log |V|} P(b_l | \Phi(v_i)) = \prod_{l=1}^{\log |V|} 1 / (1 + e^{-\Phi(v_i) \cdot \Psi(b_l)})$$



# Parameter Learning

- Randomly initialize the representations
- Each classifier in the hierarchy has a set of weights
- Use SGD (stochastic gradient descent) to update both classifier weights and vertex representations simultaneously

$$\mathcal{L} = \sum_{v \in V} \sum_{c \in W_v} -\log(P(c|v))$$

$$p(c|v) = \frac{\exp(\mathbf{z}_v^\top \mathbf{z}_c)}{\sum_{u \in V} \exp(\mathbf{z}_v^\top \mathbf{z}_u)}$$

# Results: BlogCatalog

Name	BLOGCATALOG
$ V $	10,312
$ E $	333,983
$ \mathcal{Y} $	39
	Labels
	Interests

	% Labeled Nodes	10%	20%	30%	40%	50%	60%	70%	80%	90%
Micro-F1(%)	DEEPWALK	<b>36.00</b>	<b>38.20</b>	<b>39.60</b>	<b>40.30</b>	<b>41.00</b>	<b>41.30</b>	41.50	41.50	42.00
	SpectralClustering	31.06	34.95	37.27	38.93	39.97	40.99	<b>41.66</b>	<b>42.42</b>	<b>42.62</b>
	EdgeCluster	27.94	30.76	31.85	32.99	34.12	35.00	34.63	35.99	36.29
	Modularity	27.35	30.74	31.77	32.97	34.09	36.13	36.08	37.23	38.18
	wvRN	19.51	24.34	25.62	28.82	30.37	31.81	32.19	33.33	34.28
	Majority	16.51	16.66	16.61	16.70	16.91	16.99	16.92	16.49	17.26
Macro-F1(%)	DEEPWALK	<b>21.30</b>	<b>23.80</b>	25.30	26.30	27.30	27.60	27.90	28.20	28.90
	SpectralClustering	19.14	23.57	<b>25.97</b>	<b>27.46</b>	<b>28.31</b>	<b>29.46</b>	<b>30.13</b>	<b>31.38</b>	<b>31.78</b>
	EdgeCluster	16.16	19.16	20.48	22.00	23.00	23.64	23.82	24.61	24.92
	Modularity	17.36	20.00	20.80	21.85	22.65	23.41	23.89	24.20	24.97
	wvRN	6.25	10.13	11.64	14.24	15.86	17.18	17.98	18.86	19.57
	Majority	2.52	2.55	2.52	2.58	2.58	2.63	2.61	2.48	2.62

- Feed the learned representation for node classification
- DeepWalk (node representation learning) **performs well**, especially when **labels are sparse**

# Results: YouTube

Name	YOUTUBE
$ V $	1,138,499
$ E $	2,990,443
$ \mathcal{Y} $	47
Labels	Groups

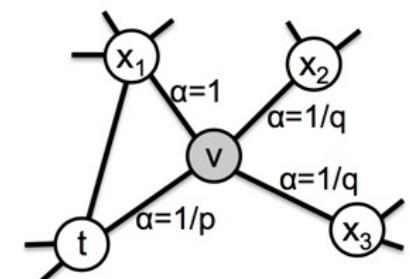
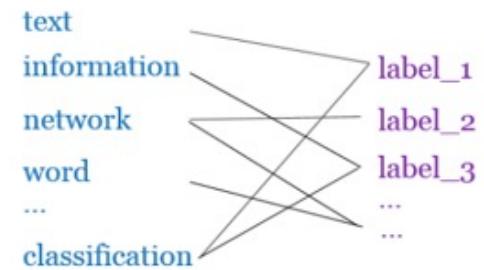
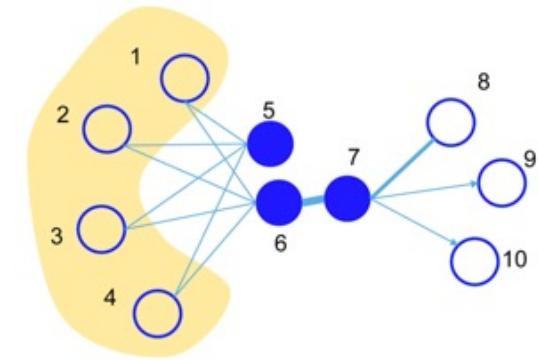
	% Labeled Nodes	1%	2%	3%	4%	5%	6%	7%	8%	9%	10%
Micro-F1(%)	DEEPWALK	<b>37.95</b>	<b>39.28</b>	<b>40.08</b>	<b>40.78</b>	<b>41.32</b>	<b>41.72</b>	<b>42.12</b>	<b>42.48</b>	<b>42.78</b>	<b>43.05</b>
	SpectralClustering	—	—	—	—	—	—	—	—	—	—
	EdgeCluster	23.90	31.68	35.53	36.76	37.81	38.63	38.94	39.46	39.92	40.07
	Modularity	—	—	—	—	—	—	—	—	—	—
	wvRN	26.79	29.18	33.1	32.88	35.76	37.38	38.21	37.75	38.68	39.42
	Majority	24.90	24.84	25.25	25.23	25.22	25.33	25.31	25.34	25.38	25.38
Macro-F1(%)	DEEPWALK	<b>29.22</b>	<b>31.83</b>	<b>33.06</b>	<b>33.90</b>	<b>34.35</b>	<b>34.66</b>	<b>34.96</b>	<b>35.22</b>	<b>35.42</b>	<b>35.67</b>
	SpectralClustering	—	—	—	—	—	—	—	—	—	—
	EdgeCluster	19.48	25.01	28.15	29.17	29.82	30.65	30.75	31.23	31.45	31.54
	Modularity	—	—	—	—	—	—	—	—	—	—
	wvRN	13.15	15.78	19.66	20.9	23.31	25.43	27.08	26.48	28.33	28.89
	Majority	6.12	5.86	6.21	6.1	6.07	6.19	6.17	6.16	6.18	6.19

- Similar results on YouTube
- Spectral Clustering does not scale to large graphs

- DeepWalk utilizes fixed-length, unbiased random walks to generate context for each node, can we do better?

# Later...

- LINE<sup>[1]</sup>: explicitly preserves both *first-order* and *second-order* proximities.
- PTE<sup>[2]</sup>: learn **heterogeneous** text network embedding via a semi-supervised manner.
- Node2vec<sup>[3]</sup>: use a **biased** random walk to better explore node's neighborhood.
- Metapath2vec<sup>[4]</sup>: **meta-path-based** random walks for heterogeneous networks.
- GATNE<sup>[5]</sup>: **inductive learning** for hetero-geneous networks.



1. J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei. 2015. Line: Large-scale information network embedding. *WWW'15*, 1067–1077.
2. J. Tang, M. Qu, and Q. Mei. 2015. Pte: Predictive text embedding through large-scale heterogeneous text networks. *KDD'15*, 1165–1174.
3. A. Grover and J. Leskovec. 2016. node2vec: Scalable feature learning for networks. *KDD'16*, 855–864.
4. Y. Dong, C. V. Nitesh and S. Ananthram. 2017. metapath2vec: Scalable Representation Learning for Heterogeneous Networks. *KDD'14*.
5. Y. Cen, X. Zou, J. Zhang, H. Yang, J. Zhou, and J. Tang. Representation Learning for Attributed Multiplex Heterogeneous Network. *KDD'19*.

# LINE: First-order proximity

## Definition

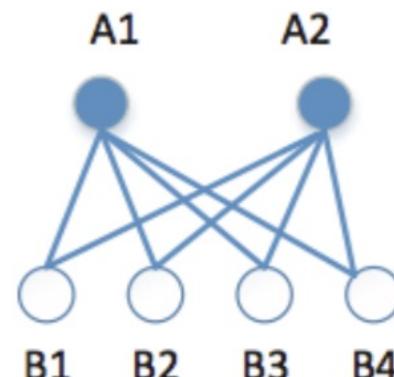
- The **first-order** proximity in a network is the **local** pairwise proximity between two vertices.
- For each pair of vertices linked by an edge  $(v_i, v_j)$ ,  $w_{ij}$  is the weight of the edge: indicates the first-order proximity between  $v_i$  and  $v_j$ .
- If no edge is observed between  $v_i$  and  $v_j$ , their **first-order** proximity is 0.



# LINE: Second-order proximity

## Definition

- The **second-order** proximity between a pair of vertices  $(v_i, v_j)$  in a network is the similarity between their neighborhood network structures.
- If no vertex is linked from/to both  $v_i$  and  $v_j$ , the **second-order** proximity between  $v_i$  and  $v_j$  is 0.



# LINE: Information Network Embedding

- Given a large network  $G = (V, E)$ , LINE aims to represent each vertex  $v \in V$  into a low-dimensional space  $R^d$
- Goal: learning a function  $f_G : V \rightarrow R^d$ ,  $d \ll |V|$ .  
In  $R^d$ , both the first-order proximity and the second-order proximity between the vertices are preserved
- For each node  $v_i \in V$ , we use  $\vec{u_i} \in R^d$  to represent the corresponding low-dimensional vector representation.

# LINE with First-order Proximity

- For each undirected  $e_{ij}$ , we define the **joint probability** between vertex  $v_i$  and  $v_j$  as follows:

$$p_1(v_i, v_j) = \frac{1}{1 + \exp(-\vec{u}_i^T \cdot \vec{u}_j)}$$

- It defines distribution  $p(\cdot, \cdot)$  over the space  $V \times V$ , and its empirical probability can be defined as  $\hat{p}_1(i, j) = \frac{w_{i,j}}{W}$  where  $W = \sum_{i,j \in E} w_{i,j}$
- Objective Function:** Minimize the following **objective function**, where  $d(\cdot, \cdot)$  is the distance between two distributions.

$$O_1 = d(\hat{p}_1(\cdot, \cdot), p_1(\cdot, \cdot))$$

- Instantiate  $d(\cdot, \cdot)$  as KL-divergence:

$$O_1 = - \sum_{(i,j) \in E} w_{ij} \log p_1(v_i, v_j)$$

# LINE with Second-order Proximity

- We first define the probability of “context”  $v_j$  generated by vertex  $v_i$  as:

$$p_2(v_j \mid v_i) = \frac{\exp(\vec{u'}_j^T \cdot \vec{u}_i)}{\sum_{k=1}^{|V|} \exp(\vec{u'}_k^T \cdot \vec{u}_i)}$$

- We minimize the following objective function:

$$O_2 = \sum_{i \in V} \lambda_i d(\hat{p}_2(\cdot \mid v_i), p_2(\cdot \mid v_i))$$

The empirical probability  $\hat{p}_2(v_j \mid v_i)$  defined as:

$$\hat{p}_2(v_j \mid v_i) = \frac{w_{ij}}{d_i} \text{ where } d_i \text{ is the out-degree of vertex } i.$$

- Replacing  $d(\cdot, \cdot)$  with KL-divergence, setting  $\lambda_i = d_i$ , we have:

$$O_2 = - \sum_{(i,j) \in E} w_{ij} \log p_2(v_j \mid v_i)$$

# Model Optimization

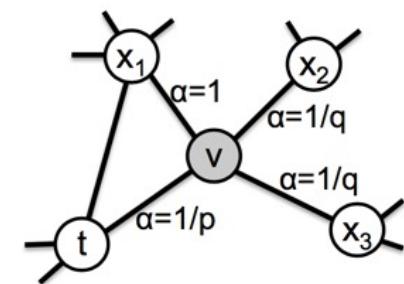
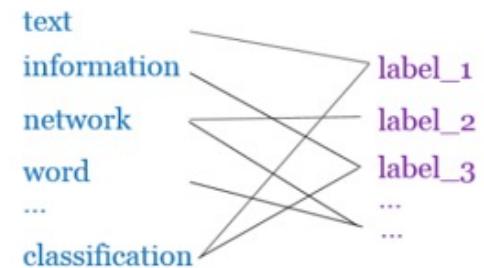
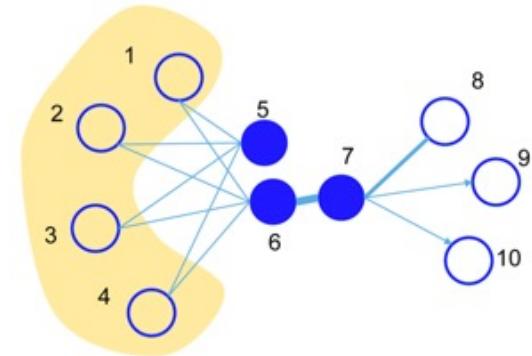
- Optimizing objective  $O_2$  is computationally expensive.
- Adopt the approach of **negative sampling**. More specifically, it specifies the following objective function for each edge  $e_{i,j}$ ,

$$\underbrace{\log\sigma(\vec{u}_j^T \cdot \vec{u}_i)}_{\text{the observed edges}} + \underbrace{\sum_{i=1}^K E_{v_n} \sim P_n(v) [\log\sigma(-\vec{u}_n^T \cdot \vec{u}_i)]}_{\text{the negative edges drawn from the noise distribution}}$$

- We can use **asynchronous stochastic gradient algorithm (ASGD)** for optimizing above Eqn.

# Later...

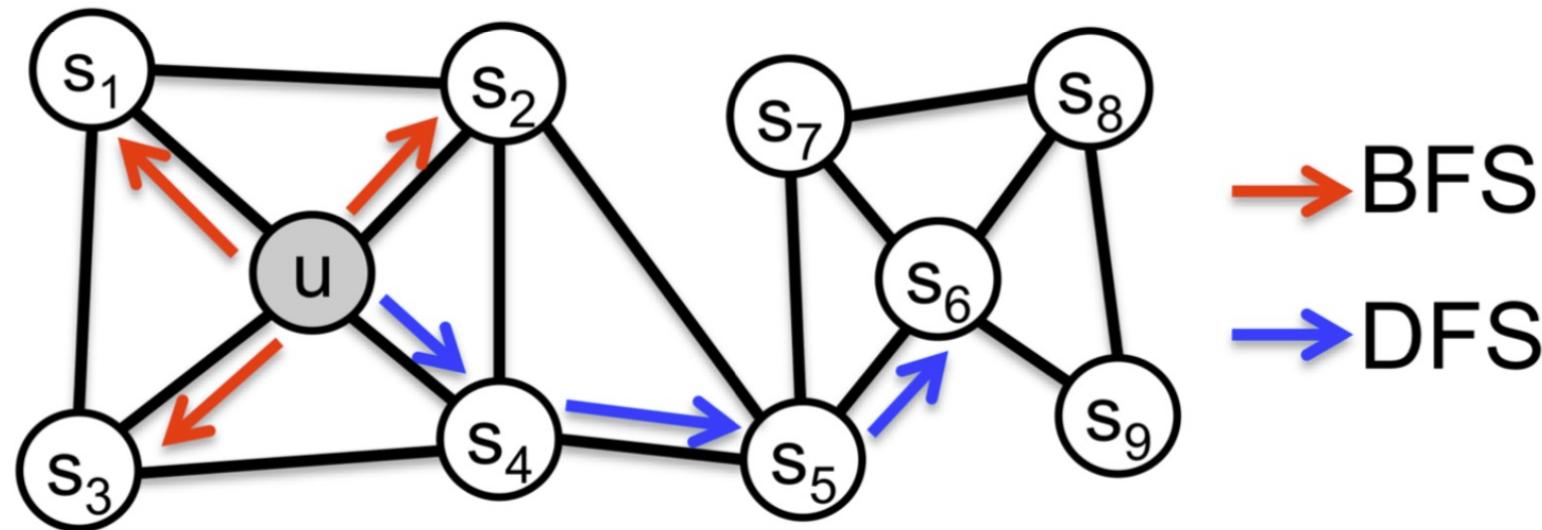
- LINE<sup>[1]</sup>: explicitly preserves both *first-order* and *second-order* proximities.
- PTE<sup>[2]</sup>: learn heterogeneous text network embedding via a semi-supervised manner.
- Node2vec<sup>[3]</sup>: use a **biased random walk** to better explore node's neighborhood.
- Metapath2vec<sup>[4]</sup>: **meta-path-based** random walks for heterogeneous networks.
- GATNE<sup>[5]</sup>: **inductive learning** for hetero-geneous networks.



1. J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei. 2015. Line: Large-scale information network embedding. *WWW'15*, 1067–1077.
2. J. Tang, M. Qu, and Q. Mei. 2015. Pte: Predictive text embedding through large-scale heterogeneous text networks. *KDD'15*, 1165–1174.
3. A. Grover and J. Leskovec. 2016. node2vec: Scalable feature learning for networks. *KDD'16*, 855–864.
4. Y. Dong, C. V. Nitesh and S. Ananthram. 2017. metapath2vec: Scalable Representation Learning for Heterogeneous Networks. *KDD'14*.
5. Y. Cen, X. Zou, J. Zhang, H. Yang, J. Zhou, and J. Tang. Representation Learning for Attributed Multiplex Heterogeneous Network. *KDD'19*.

# node2vec: Biased Walks

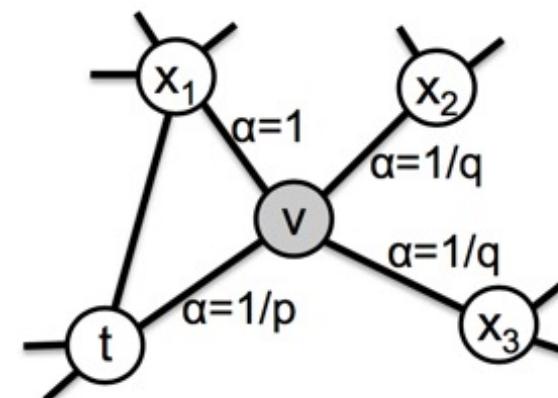
- Use biased random walks to trade off **local and global** views of the network



- Biased walks is a special case of random walk, thus node2vec is a special case of DeepWalk

# node2vec

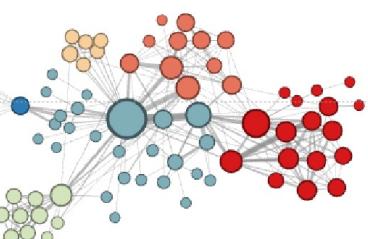
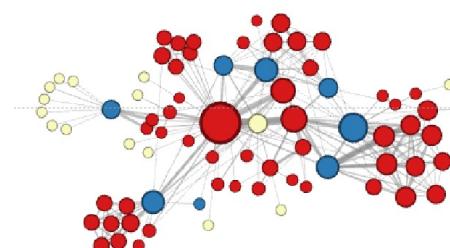
- Biased random walk  $R$  that given a node  $v$  generates random walk neighborhood  $N_{rw}(v)$
- Return parameter  $p$ :
  - Return back to the previous node
- In-out parameter  $q$ :
  - Moving outwards (DFS) vs. inwards (BFS)



# node2vec

- Biased random walk  $R$  that given a node  $v$  generates random walk neighborhood  $N_{rw}(v)$
- Return parameter  $p$ :
  - Return back to the previous node
- In-out parameter  $q$ :
  - Moving outwards (DFS) vs. inwards (BFS)

Interactions of characters in a novel:



$p=1, q=2$

Microscopic view of the  
network neighbourhood

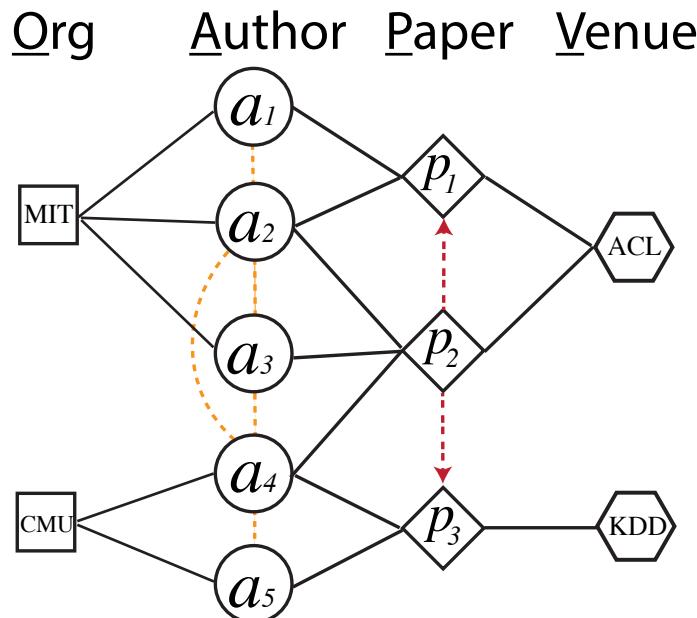
$p=1, q=0.5$

Macroscopic view of the  
network neighbourhood

Picture snipped from Leskovec

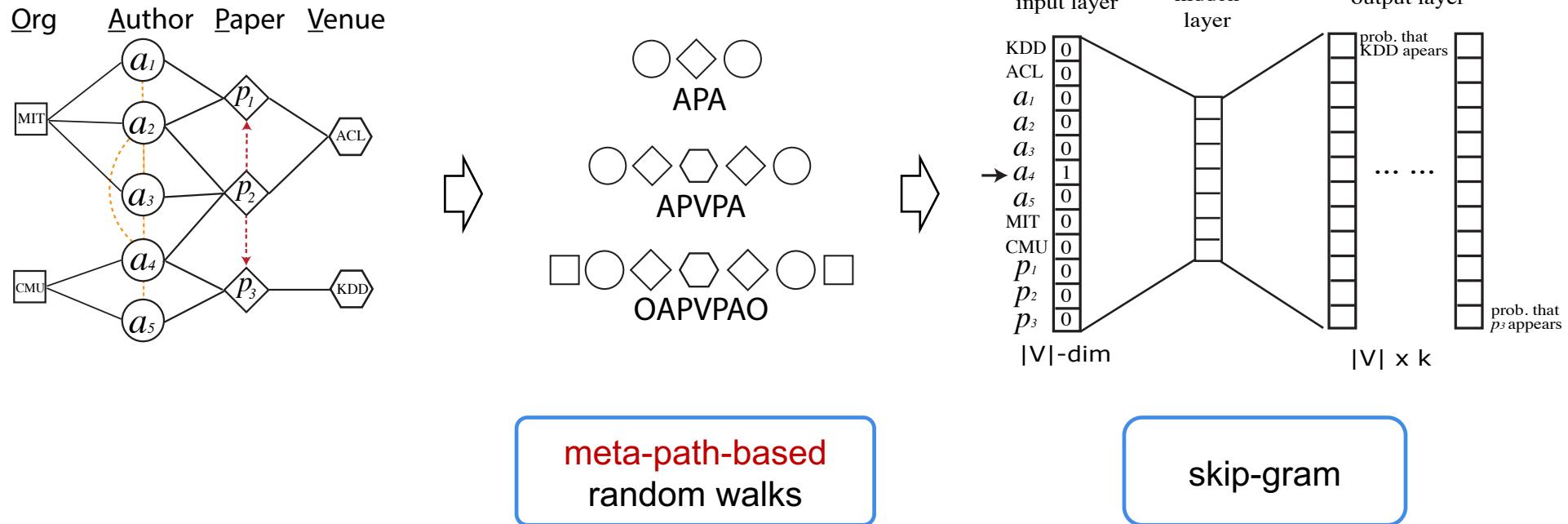
# Metapath2vec: Heterogeneous Random Walk

- Input: a heterogeneous graph  $G = (V, E)$
- Output:  $X \in R^{|V| \times k}$ ,  $k \ll |V|$ ,  $k$ -dim vector  $X_v$  for each node  $v$ .



- How do we random walk over **heterogeneous** networks?
- How do we apply skip-gram over **different types** of nodes?

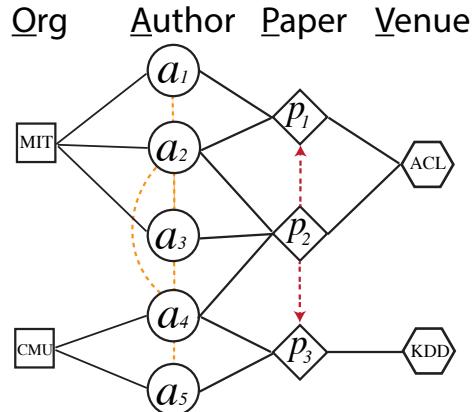
# Metapath2vec: Heterogeneous Random Walk



1. Sun and Han. Mining heterogeneous information networks: Principles and Methodologies. Morgan & Claypool Publishers, 2012.
2. Dong et al. metapath2vec: scalable representation learning for heterogeneous networks. In *ACM KDD 2017*. The most cited paper in KDD'17 as of 2018.

# Metapath2vec: Heterogeneous Random Walk

- Given a meta-path scheme



$$\mathcal{P}: V_1 \xrightarrow{R_1} V_2 \xrightarrow{R_2} \cdots V_t \xrightarrow{R_t} V_{t+1} \cdots \xrightarrow{R_{l-1}} V_l$$

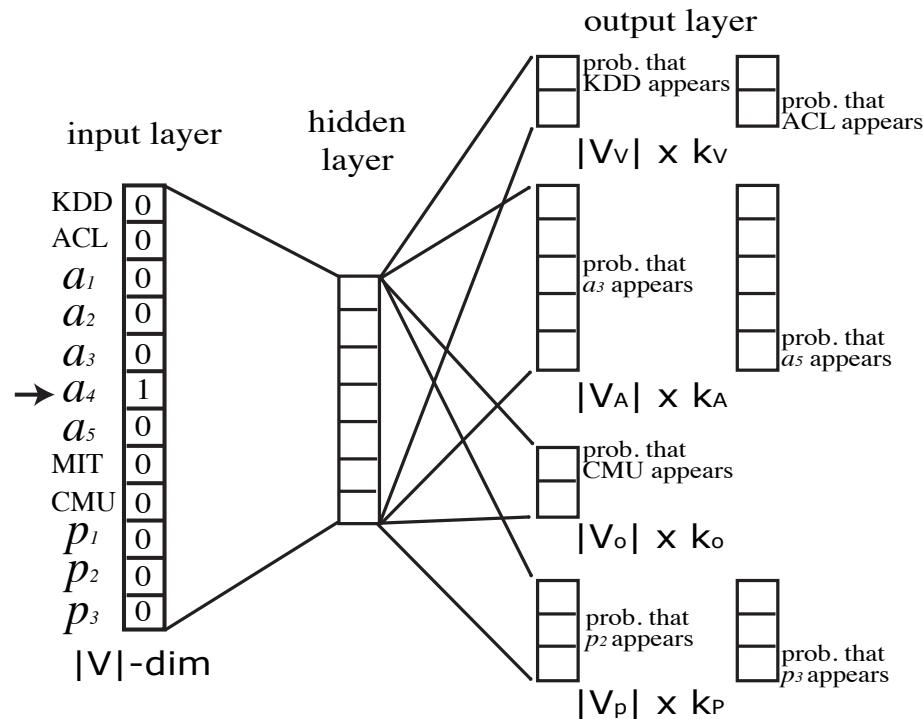
- The transition probability at step  $i$  is defined as

$$p(v^{i+1}|v_t^i, \mathcal{P}) = \begin{cases} \frac{1}{|N_{t+1}(v_t^i)|} & (v^{i+1}, v_t^i) \in E, \phi(v^{i+1}) = t+1 \\ 0 & (v^{i+1}, v_t^i) \in E, \phi(v^{i+1}) \neq t+1 \\ 0 & (v^{i+1}, v_t^i) \notin E \end{cases}$$

- Recursive guidance for random walkers, i.e.,

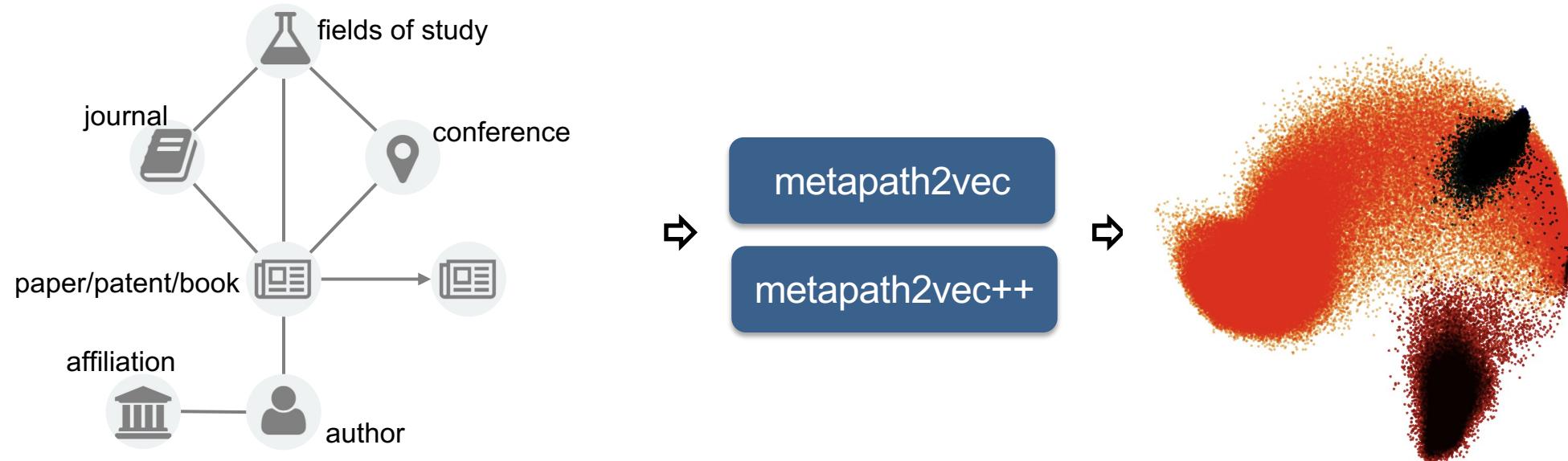
$$p(v^{i+1}|v_t^i) = p(v^{i+1}|v_1^i), \text{ if } t = l$$

# Metapath2vec: Heterogeneous Random Walk



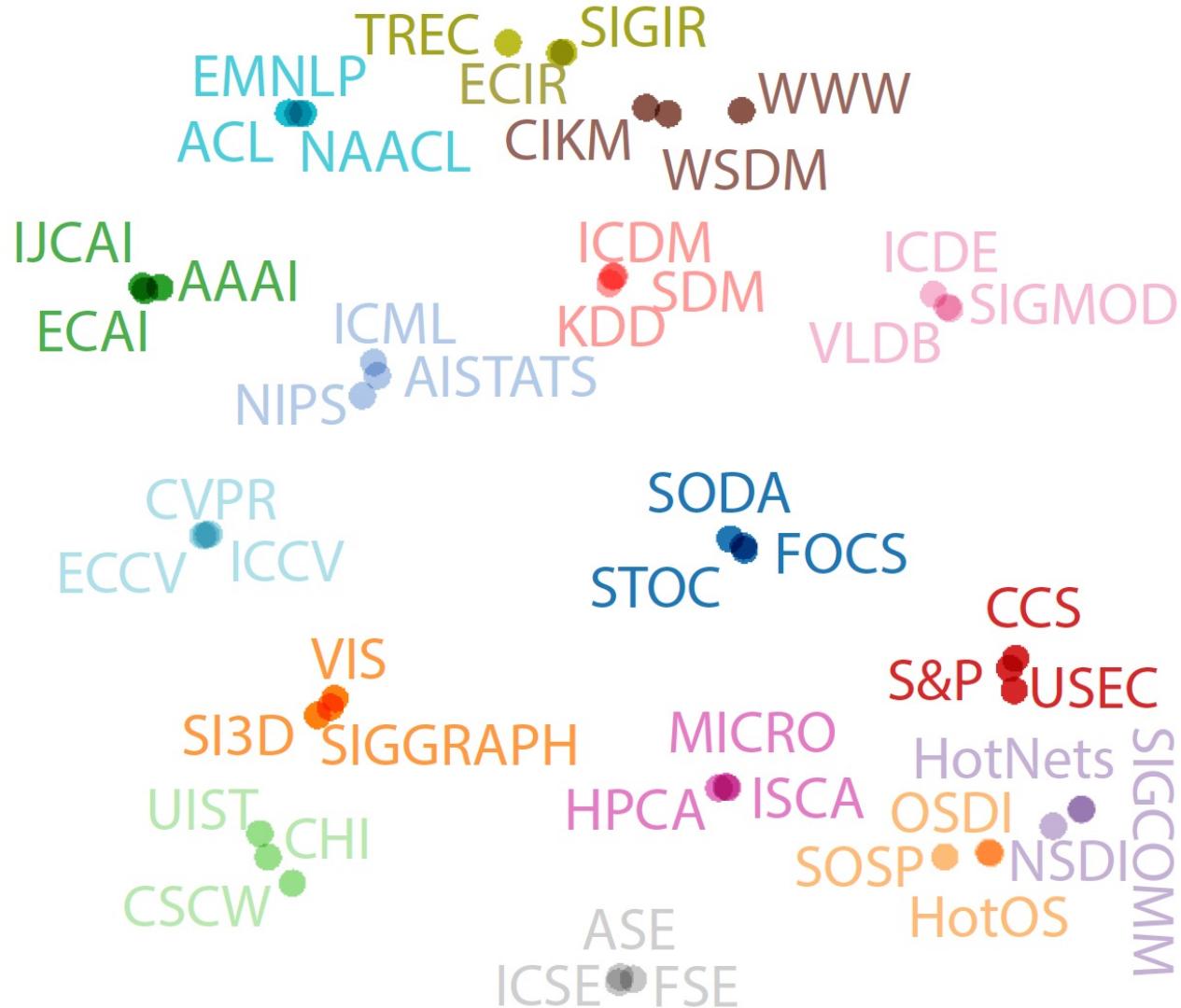
$$\mathcal{O}(\mathbf{X}) = \log \sigma(X_{c_t} \cdot X_v) + \sum_{k=1}^K \mathbb{E}_{u_t^k \sim P_t(u_t)} [\log \sigma(-X_{u_t^k} \cdot X_v)]$$

# Application: Embedding Academic Graph

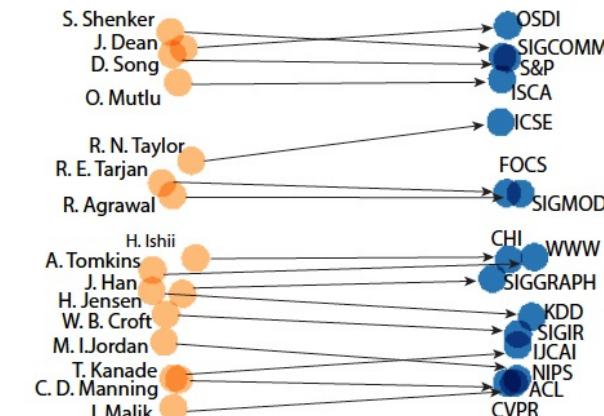
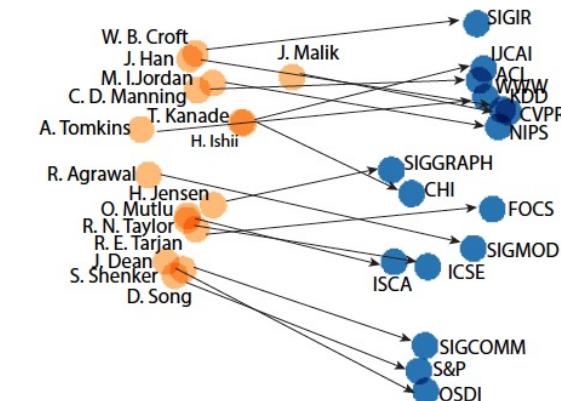
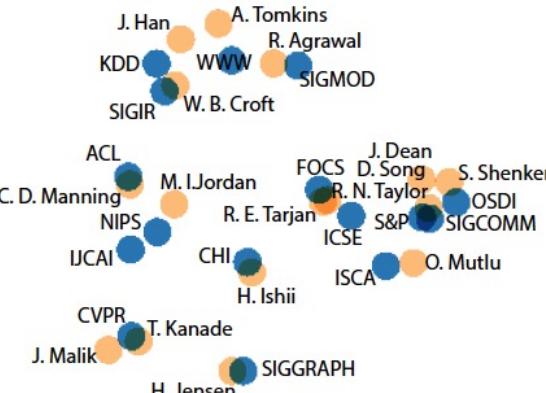
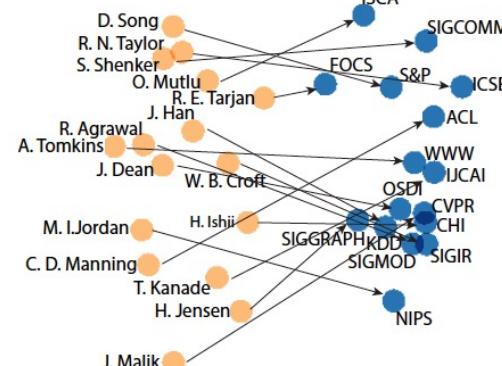
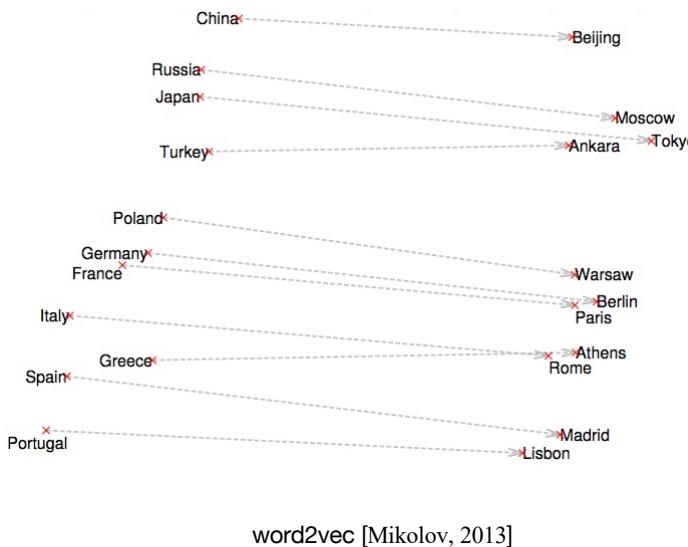


Microsoft Academic Graph  
&  
AMiner

# Application 2: Node Clustering

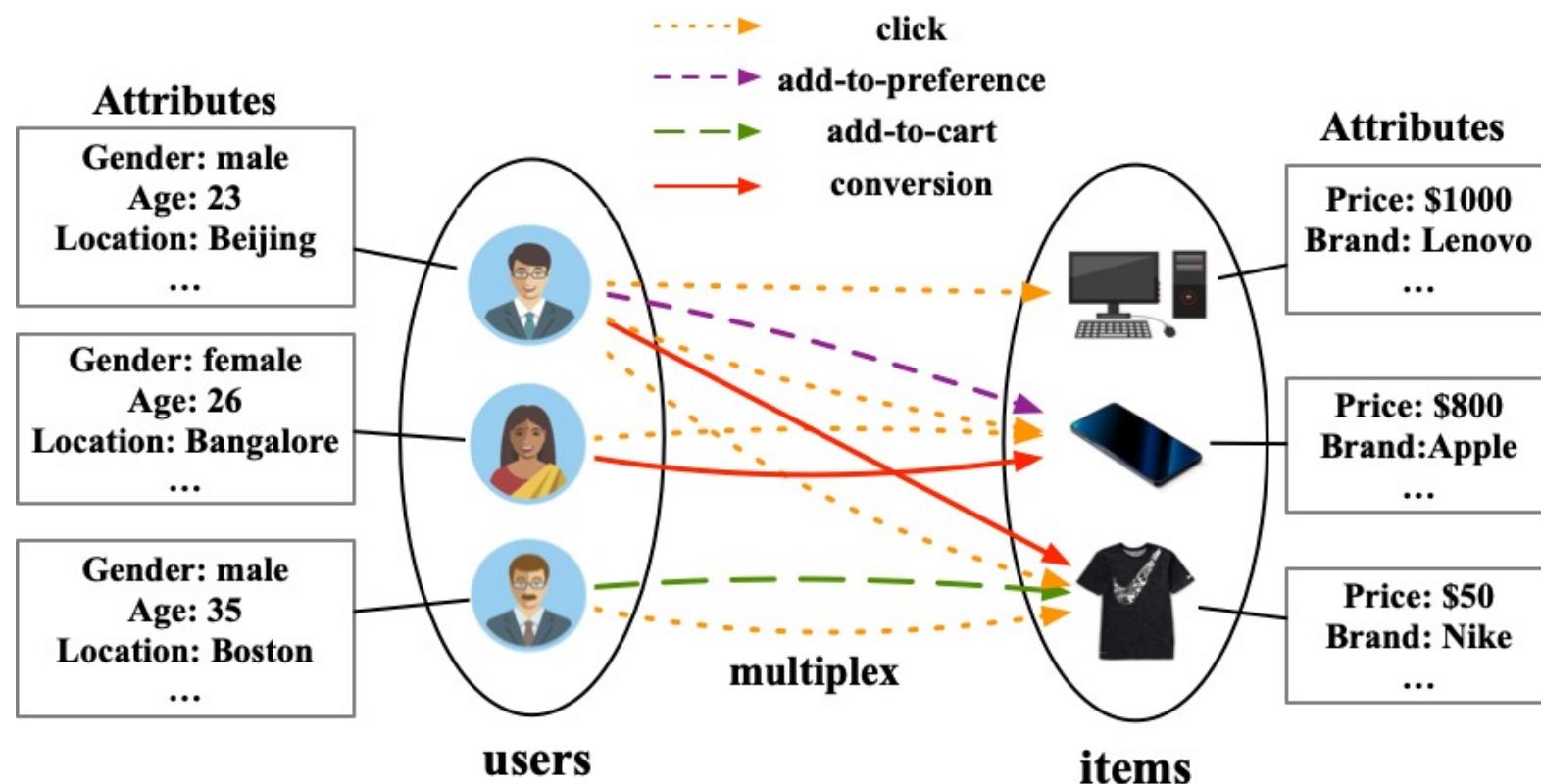


# Visualization



# GATNE: Attributed Multiplex Heterogeneous Network Embedding

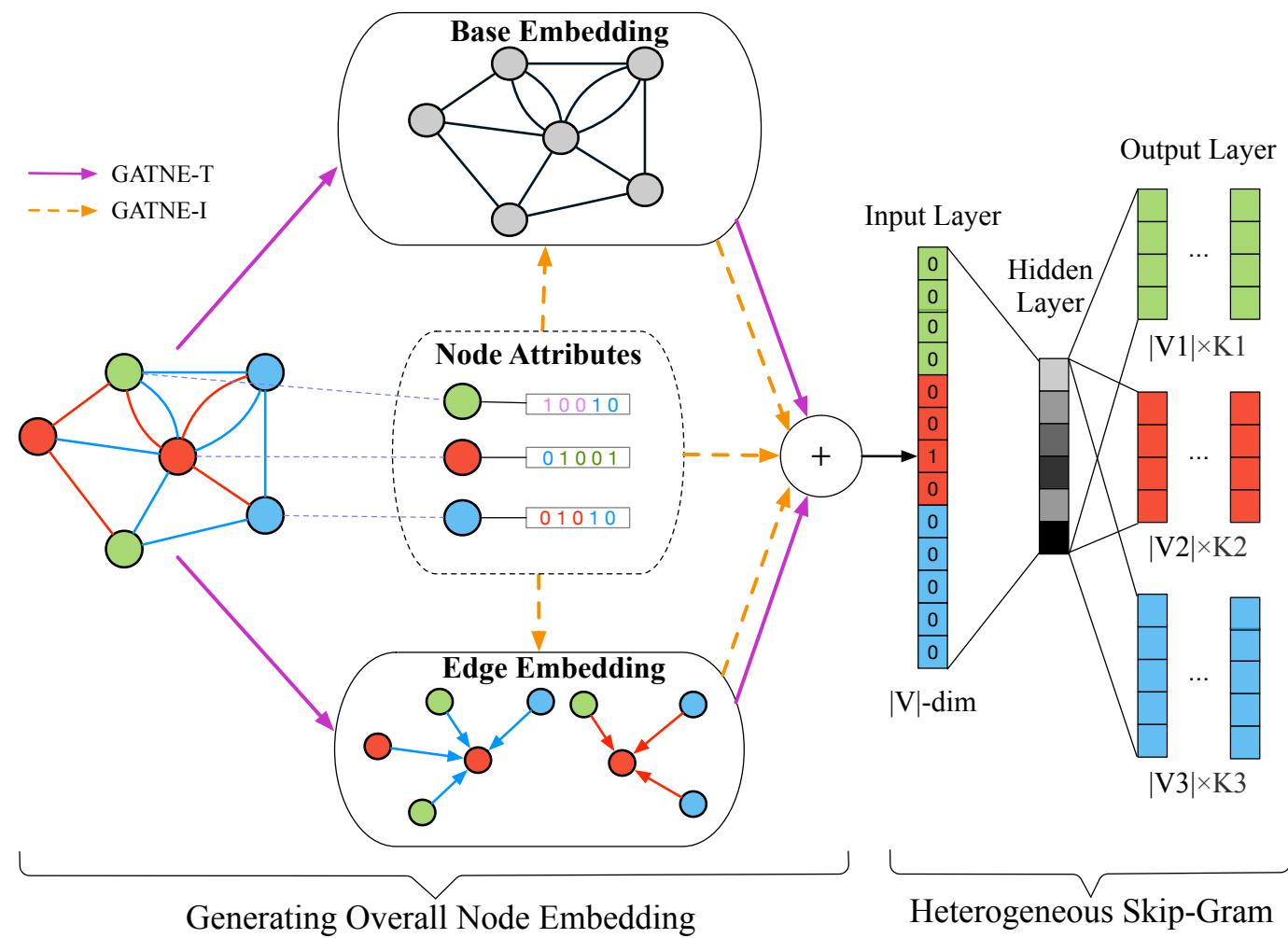
- Recommend items to users by considering *Attributed Multiplex Heterogeneous Networks (AMHEN)*



# Different Types of Network Embedding

Network Type	Method	Heterogeneity		Attribute
		Node Type	Edge Type	
Homogeneous Network (HON)	DeepWalk [27] LINE [35] node2vec [10] NetMF [29] NetSMF [28]	Single	Single	/
Attributed Homogeneous Network (AHON)	TADW [41] LANE [16] AANE [15] SNE [20] DANE [9] ANRL [44]	Single	Single	Attributed
Heterogeneous Network (HEN)	PTE [34] metapath2vec [7] HERec [31]	Multi	Single	/
Attributed HEN (AHEN)	HNE [3]	Multi	Single	Attributed
Multiplex Heterogeneous Network (MHEN)	PMNE [22] MVE [30] MNE [43] mvn2vec [32]	Single	Multi	/
	GATNE-T	Multi	Multi	
Attributed MHEN (AMHEN)	GATNE-I	Multi	Multi	Attributed

# Multiplex Heterogeneous Graph Embedding



# Recommendation Results

- Data
  - Small: Amazon, YouTube, Twitter w/ 10K nodes
  - Large: Alibaba w/ **40M nodes** and **0.5B edges**

	Amazon			YouTube			Twitter			Alibaba-S		
	ROC-AUC	PR-AUC	F1									
DeepWalk	94.20	94.03	87.38	71.11	70.04	65.52	69.42	72.58	62.68	59.39	60.62	56.10
node2vec	94.47	94.30	87.88	71.21	70.32	65.36	69.90	73.04	63.12	62.26	63.40	58.49
LINE	81.45	74.97	76.35	64.24	63.25	62.35	62.29	60.88	58.18	53.97	54.65	52.85
metapath2vec	94.15	94.01	87.48	70.98	70.02	65.34	69.35	72.61	62.70	60.94	61.40	58.25
ANRL	71.68	70.30	67.72	75.93	73.21	70.65	70.04	67.16	64.69	58.17	55.94	56.22
PMNE(n)	95.59	95.48	89.37	65.06	63.59	60.85	69.48	72.66	62.88	62.23	63.35	58.74
PMNE(r)	88.38	88.56	79.67	70.61	69.82	65.39	62.91	67.85	56.13	55.29	57.49	53.65
PMNE(c)	93.55	93.46	86.42	68.63	68.22	63.54	67.04	70.23	60.84	51.57	51.78	51.44
MVE	92.98	93.05	87.80	70.39	70.10	65.10	72.62	73.47	67.04	60.24	60.51	57.08
MNE	90.28	91.74	83.25	82.30	82.18	75.03	91.37	91.65	84.32	62.79	63.82	58.74
GATNE-T	<b>97.44</b>	<b>97.05</b>	<b>92.87</b>	<b>84.61</b>	81.93	<b>76.83</b>	<b>92.30</b>	91.77	<b>84.96</b>	66.71	67.55	62.48
GATNE-I	96.25	94.77	91.36	84.47	<b>82.32</b>	<b>76.83</b>	92.04	<b>91.95</b>	84.38	<b>70.87</b>	<b>71.65</b>	<b>65.54</b>

GATNE outperforms all sorts of baselines in the various datasets.

\*\* Code available at <https://github.com/THUDM/GATNE>

# Alibaba Offline A/B Tests

- GATNE-I is deployed on Alibaba's distributed cloud platform for its recommendation system. The training dataset has about 100 million users and 10 million items with 10 billion interactions between them.
- Under the framework of A/B tests, an offline test is conducted on GATNE-I, MNE and DeepWalk. The experimental goal is to maximize Hit-Rate. The results demonstrate that GATNE-I improves Hit-Rate by **3.26% and 24.26%** compared to MNE and DeepWalk respectively.

# Questions

- What are the **fundamentals** underlying the different models?
- Can we **unify** the different graph embedding approaches?

# Unifying DeepWalk, LINE, PTE, and node2vec into Matrix Forms

Algorithm	Closed Matrix Form
DeepWalk	$\log \left( \text{vol}(G) \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right) - \log b$
LINE	$\log (\text{vol}(G) \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1}) - \log b$
PTE	$\log \begin{pmatrix} \alpha \text{vol}(G_{\text{ww}}) (\mathbf{D}_{\text{row}}^{\text{ww}})^{-1} \mathbf{A}_{\text{ww}} (\mathbf{D}_{\text{col}}^{\text{ww}})^{-1} \\ \beta \text{vol}(G_{\text{dw}}) (\mathbf{D}_{\text{row}}^{\text{dw}})^{-1} \mathbf{A}_{\text{dw}} (\mathbf{D}_{\text{col}}^{\text{dw}})^{-1} \\ \gamma \text{vol}(G_{\text{lw}}) (\mathbf{D}_{\text{row}}^{\text{lw}})^{-1} \mathbf{A}_{\text{lw}} (\mathbf{D}_{\text{col}}^{\text{lw}})^{-1} \end{pmatrix} - \log b$
node2vec	$\log \left( \frac{\frac{1}{2T} \sum_{r=1}^T \left( \sum_u \mathbf{X}_{w,u} \underline{\mathbf{P}}_{c,w,u}^r + \sum_u \mathbf{X}_{c,u} \underline{\mathbf{P}}_{w,c,u}^r \right)}{(\sum_u \mathbf{X}_{w,u})(\sum_u \mathbf{X}_{c,u})} \right) - \log b$

$\mathbf{A}$ :  $\mathbf{A} \in \mathbb{R}_+^{|V| \times |V|}$  is  $G$ 's adjacency matrix with  $A_{i,j}$  as the edge weight between vertices  $i$  and  $j$ ;

$D_{\text{col}}$ :  $D_{\text{col}} = \text{diag}(\mathbf{A}^\top \mathbf{e})$  is the diagonal matrix with column sum of  $\mathbf{A}$ ;

$D_{\text{row}}$ :  $D_{\text{row}} = \text{diag}(\mathbf{A}\mathbf{e})$  is the diagonal matrix with row sum of  $\mathbf{A}$ ;

$\mathbf{D}$ : For undirected graphs ( $\mathbf{A}^\top = \mathbf{A}$ ),  $D_{\text{col}} = D_{\text{row}}$ . For brevity,  $\mathbf{D}$  represents both  $D_{\text{col}}$  &  $D_{\text{row}}$ .

$\mathbf{D} = \text{diag}(d_1, \dots, d_{|V|})$ , where  $d_i$  represents generalized degree of vertex  $i$ ;

$\text{vol}(G)$ :  $\text{vol}(G) = \sum_i \sum_j A_{i,j} = \sum_i d_i$  is the volume of an weighted graph  $G$ ;

$T$  &  $b$ : The context window size and the number of negative sampling in skip-gram, respectively.

# DeepWalk is factorizing a matrix



DeepWalk is asymptotically and implicitly factorizing

$$\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

$$\text{vol}(G) = \sum_i \sum_j A_{ij}$$

$\mathbf{A}$  Adjacency matrix

$b$ : #negative samples

$\mathbf{D}$  Degree matrix

$T$ : context window size

# Skip gram with negative sampling

## Skip-gram with negative sampling (SGNS)

- SGNS maintains a multiset  $\mathcal{D}$  that counts the occurrence of each word-context pair  $(w, c)$
- Objective

$$\mathcal{L} = \sum_w \sum_c (\#(w, c) \log g(x_w^T x_c) + \frac{b\#(w)\#(c)}{|\mathcal{D}|} \log g(-x_w^T x_c))$$

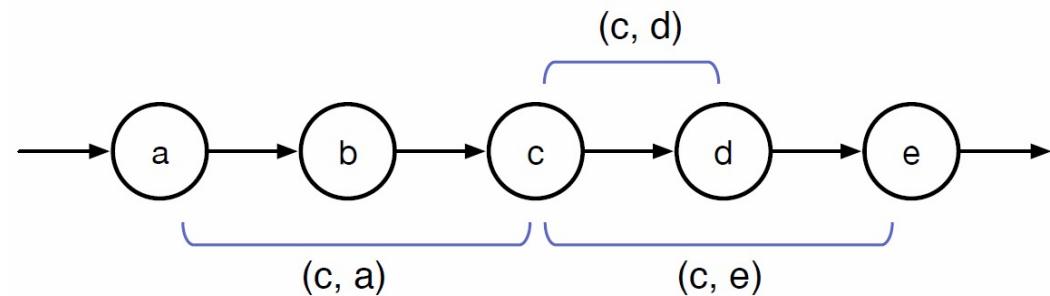
- For sufficiently large dimension  $d$ , the objective above is equivalent to factorizing the PMI matrix

$$\log \frac{\#(w, c)|\mathcal{D}|}{b\#(w)\#(c)}$$

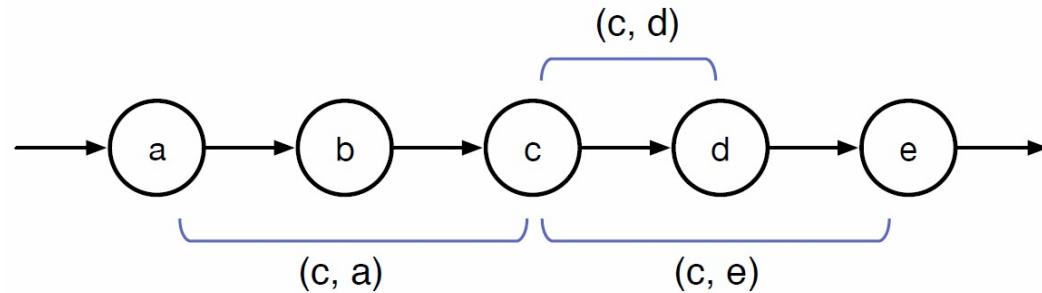
# Understanding random walk + skip gram

- 1 **for**  $n = 1, 2, \dots, N$  **do**
- 2     | Pick  $w_1^n$  according to a probability distribution  $P(w_1)$ ;
- 3     | Generate a vertex sequence  $(w_1^n, \dots, w_L^n)$  of length  $L$  by a random walk on network  $G$ ;
- 4     | **for**  $j = 1, 2, \dots, L - T$  **do**
- 5         | | **for**  $r = 1, \dots, T$  **do**
- 6             | | | Add vertex-context pair  $(w_j^n, w_{j+r}^n)$  to multiset  $\mathcal{D}$ ;
- 7             | | | Add vertex-context pair  $(w_{j+r}^n, w_j^n)$  to multiset  $\mathcal{D}$ ;
- 8     | | | | Run SGNS on  $\mathcal{D}$  with  $b$  negative samples.

# Understanding random walk + skip gram

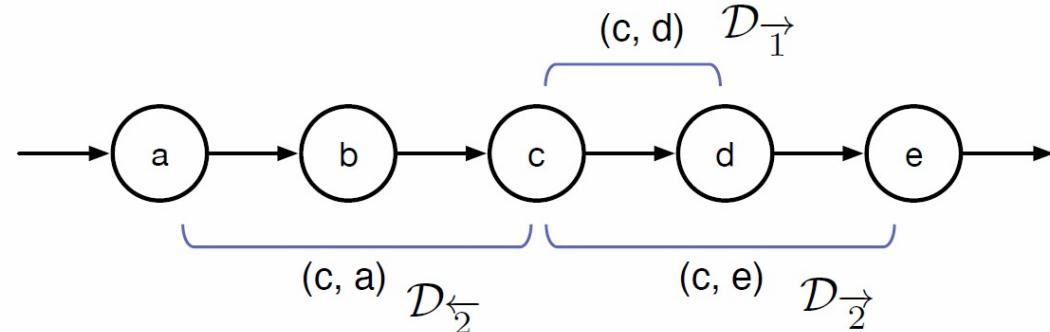


# Understanding random walk + skip gram



Suppose the multiset  $\mathcal{D}$  is constructed based on random walk on graphs, can we interpret  $\log \frac{\#(w,c)|\mathcal{D}|}{b\#(w)\#(c)}$  with graph structures?

# Understanding random walk + skip gram



- Partition the multiset  $\mathcal{D}$  into several sub-multisets according to the way in which each node and its context appear in a random walk node sequence.
- More formally, for  $r = 1, 2, \dots, T$ , we define

$$\mathcal{D}_{\rightarrow} = \{(w, c) : (w, c) \in \mathcal{D}, w = w_j^n, c = w_{j+r}^n\}$$

$$\mathcal{D}_{\leftarrow} = \{(w, c) : (w, c) \in \mathcal{D}, w = w_{j+r}^n, c = w_j^n\}$$

Distinguish direction  
and distance

# Understanding random walk + skip gram

$$\log \left( \frac{\#(w, c) |\mathcal{D}|}{b \#(w) \cdot \#(c)} \right) = \log \left( \frac{\frac{\#(w, c)}{|\mathcal{D}|}}{b \frac{\#(w)}{|\mathcal{D}|} \frac{\#(c)}{|\mathcal{D}|}} \right)$$

the length of random walk  $L \rightarrow \infty$

$$\frac{\#(w, c)}{|\mathcal{D}|} \xrightarrow{p} \frac{1}{2T} \sum_{r=1}^T \left( \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c} + \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w} \right)$$

$$\frac{\#(w)}{|\mathcal{D}|} \xrightarrow{p} \frac{d_w}{\text{vol}(G)}$$

$$\frac{\#(c)}{|\mathcal{D}|} \xrightarrow{p} \frac{d_c}{\text{vol}(G)}$$

$$\mathbf{P} = \mathbf{D}^{-1} \mathbf{A}$$

$$\begin{aligned} \frac{\#(w, c) |\mathcal{D}|}{\#(w) \cdot \#(c)} &= \frac{\frac{\#(w, c)}{|\mathcal{D}|}}{\frac{\#(w)}{|\mathcal{D}|} \cdot \frac{\#(c)}{|\mathcal{D}|}} \xrightarrow{p} \frac{\frac{1}{2T} \sum_{r=1}^T \left( \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c} + \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w} \right)}{\frac{d_w}{\text{vol}(G)} \cdot \frac{d_c}{\text{vol}(G)}} \\ &= \frac{\text{vol}(G)}{2T} \left( \frac{1}{d_c} \sum_{r=1}^T (\mathbf{P}^r)_{w,c} + \frac{1}{d_w} \sum_{r=1}^T (\mathbf{P}^r)_{c,w} \right) \end{aligned}$$

$$\frac{\#(w,c)\left|\mathcal{D}\right|}{\#(w)\cdot \#(c)} \overset{p}{\rightarrow} \frac{\text{vol}(G)}{2T}\left(\frac{1}{d_c}\sum_{r=1}^T (\boldsymbol{P}^r)_{w,c} + \frac{1}{d_w}\sum_{r=1}^T (\boldsymbol{P}^r)_{c,w}\right)$$

$$\begin{aligned}& \frac{\text{vol}(G)}{2T}\left(\sum_{r=1}^T \boldsymbol{P}^r \boldsymbol{D}^{-1} + \sum_{r=1}^T \boldsymbol{D}^{-1} (\boldsymbol{P}^r)^{\top}\right) \\&= \frac{\text{vol}(G)}{2T}\left(\sum_{r=1}^T \underbrace{\boldsymbol{D}^{-1} \boldsymbol{A} \times \cdots \times \boldsymbol{D}^{-1} \boldsymbol{A}}_{r \text{ terms}} \boldsymbol{D}^{-1} + \sum_{r=1}^T \boldsymbol{D}^{-1} \underbrace{\boldsymbol{A} \boldsymbol{D}^{-1} \times \cdots \times \boldsymbol{A} \boldsymbol{D}^{-1}}_{r \text{ terms}}\right) \\&= \frac{\text{vol}(G)}{T} \sum_{r=1}^T \underbrace{\boldsymbol{D}^{-1} \boldsymbol{A} \times \cdots \times \boldsymbol{D}^{-1} \boldsymbol{A}}_{r \text{ terms}} \boldsymbol{D}^{-1} = \text{vol}(G) \left(\frac{1}{T} \sum_{r=1}^T \boldsymbol{P}^r\right) \boldsymbol{D}^{-1}.\end{aligned}$$

# DeepWalk is factorizing a matrix



DeepWalk is asymptotically and implicitly factorizing

$$\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

$$\text{vol}(G) = \sum_i \sum_j A_{ij}$$

$\mathbf{A}$  Adjacency matrix

$b$ : #negative samples

$\mathbf{D}$  Degree matrix

$T$ : context window size

# LINE

- ▶ Objective of LINE:

$$\mathcal{L} = \sum_{i=1}^{|V|} \sum_{j=1}^{|V|} \left( \mathbf{A}_{i,j} \log g(\mathbf{x}_i^\top \mathbf{y}_j) + \frac{bd_i d_j}{\text{vol}(G)} \log g(-\mathbf{x}_i^\top \mathbf{y}_j) \right).$$

- ▶ Align it with the Objective of SGNS:

$$\mathcal{L} = \sum_w \sum_c \left( \#(w, c) \log g(\mathbf{x}_w^\top \mathbf{y}_c) + \frac{b\#(w)\#(c)}{|\mathcal{D}|} \log g(-\mathbf{x}_w^\top \mathbf{y}_c) \right).$$

- ▶ LINE is actually factorizing

$$\log \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right)$$

- ▶ Recall DeepWalk's matrix form:

$$\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right).$$

**Observation** LINE is a special case of DeepWalk ( $T = 1$ ).

# PTE

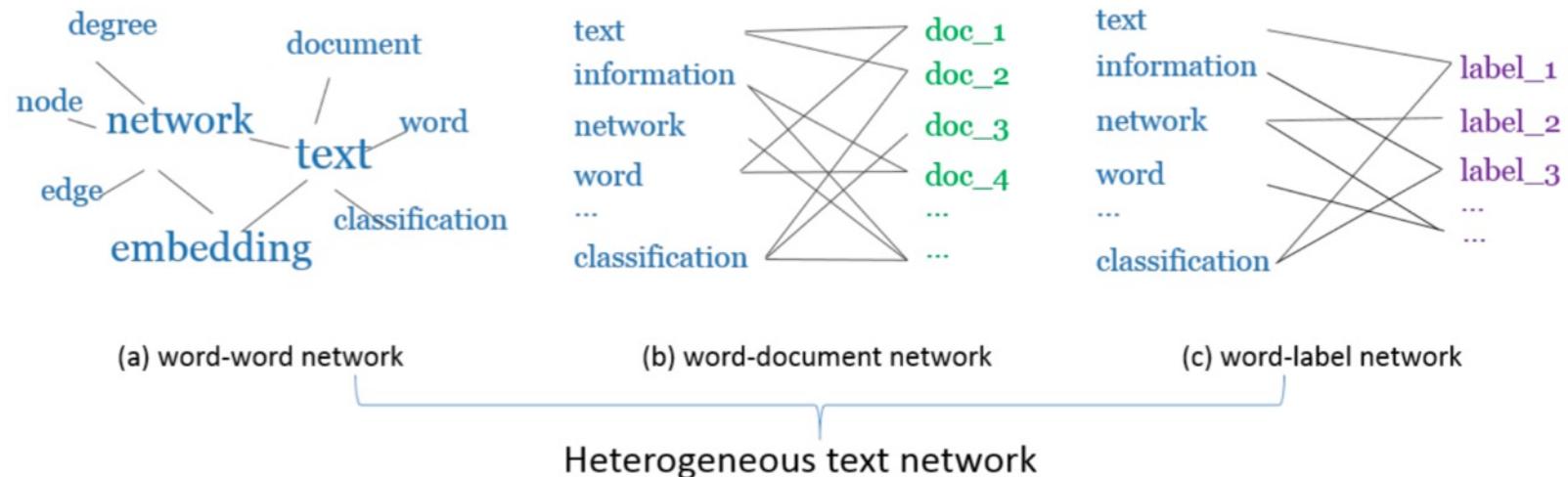


Figure 2: Heterogeneous Text Network.

$$\log \left( \begin{bmatrix} \alpha \text{vol}(G_{\text{ww}})(\mathbf{D}_{\text{row}}^{\text{ww}})^{-1} \mathbf{A}_{\text{ww}} (\mathbf{D}_{\text{col}}^{\text{ww}})^{-1} \\ \beta \text{vol}(G_{\text{dw}})(\mathbf{D}_{\text{row}}^{\text{dw}})^{-1} \mathbf{A}_{\text{dw}} (\mathbf{D}_{\text{col}}^{\text{dw}})^{-1} \\ \gamma \text{vol}(G_{\text{lw}})(\mathbf{D}_{\text{row}}^{\text{lw}})^{-1} \mathbf{A}_{\text{lw}} (\mathbf{D}_{\text{col}}^{\text{lw}})^{-1} \end{bmatrix} \right) - \log b,$$

# Understanding node2vec

$$\underline{\mathbf{T}}_{u,v,w} = \begin{cases} \frac{1}{p} & (u, v) \in E, (v, w) \in E, u = w; \\ 1 & (u, v) \in E, (v, w) \in E, u \neq w, (w, u) \in E; \\ \frac{1}{q} & (u, v) \in E, (v, w) \in E, u \neq w, (w, u) \notin E; \\ 0 & \text{otherwise.} \end{cases}$$

$$\underline{\mathbf{P}}_{u,v,w} = \text{Prob}(w_{j+1} = u | w_j = v, w_{j-1} = w) = \frac{\underline{\mathbf{T}}_{u,v,w}}{\sum_u \underline{\mathbf{T}}_{u,v,w}}.$$

## Stationary Distribution

$$\sum_w \underline{\mathbf{P}}_{u,v,w} \mathbf{X}_{v,w} = \mathbf{X}_{u,v}$$

Existence guaranteed by Perron-Frobenius theorem, but may not be unique.

# Understanding node2vec

## Theorem

*node2vec is asymptotically and implicitly factorizing a matrix whose entry at  $w$ -th row,  $c$ -th column is*

$$\log \left( \frac{\frac{1}{2T} \sum_{r=1}^T (\sum_u \mathbf{X}_{w,u} \underline{\mathbf{P}}_{c,w,u}^r + \sum_u \mathbf{X}_{c,u} \underline{\mathbf{P}}_{w,c,u}^r)}{b(\sum_u \mathbf{X}_{w,u})(\sum_u \mathbf{X}_{c,u})} \right)$$

Can we directly factorize the **derived  
matrices** for learning embeddings?

# NetMF

- DeepWalk is implicitly factorizing

$$\textcolor{red}{M} = \log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

- NetMF is explicitly factorizing

$$\textcolor{red}{M} = \log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

# NetMF

- DeepWalk is implicitly factorizing

$$\mathbf{M} = \log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

Recall that in random walk + skip-gram based embedding models:  $\mathbf{z}_v^\top \mathbf{z}_c \rightarrow$  the probability that node  $v$  and context  $c$  appear on a random walk path

- NetMF is explicitly factorizing

$$\mathbf{M} = \log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

$\mathbf{z}_v^\top \mathbf{z}_c \rightarrow$  the similarity score  $M_{vc}$  between node  $v$  and context  $c$  defined by this matrix

# NetMF

- 
- 1 Eigen-decomposition  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \approx \mathbf{U}_h \boldsymbol{\Lambda}_h \mathbf{U}_h^\top$ ;
  - 2 Approximate  $\mathbf{M}$  with
$$\hat{\mathbf{M}} = \frac{\text{vol}(G)}{b} \mathbf{D}^{-1/2} \mathbf{U}_h \left( \frac{1}{T} \sum_{r=1}^T \boldsymbol{\Lambda}_h^r \right) \mathbf{U}_h^\top \mathbf{D}^{-1/2};$$
  - 3 Compute  $\hat{\mathbf{M}}' = \max(\hat{\mathbf{M}}, 1)$ ;
  - 4 Rank- $d$  approximation by SVD:  $\log \hat{\mathbf{M}}' = \mathbf{U}_d \boldsymbol{\Sigma}_d \mathbf{V}_d^\top$ ;
  - 5 **return**  $\mathbf{U}_d \sqrt{\boldsymbol{\Sigma}_d}$  as network embedding.
- 

Approximate  $\mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$   
with its top- $h$  eigenpairs  
 $\mathbf{U}_h \boldsymbol{\Lambda}_h \mathbf{U}_h^\top$

The Arnoldi algorithm [1]  
for significant time  
reduction

# Error Bound for NetMF for a large window size $T$

- According to Frobenius norm's property

$$\begin{aligned} \left| \log M'_{i,j} - \log \hat{M}'_{i,j} \right| &= \log \frac{\hat{M}'_{i,j}}{M'_{i,j}} = \log \left( 1 + \frac{\hat{M}'_{i,j} - M'_{i,j}}{M'_{i,j}} \right) \\ &\leq \frac{\hat{M}'_{i,j} - M'_{i,j}}{M'_{i,j}} \leq \hat{M}'_{i,j} - M'_{i,j} = \left| \hat{M}'_{i,j} - M'_{i,j} \right| \end{aligned}$$

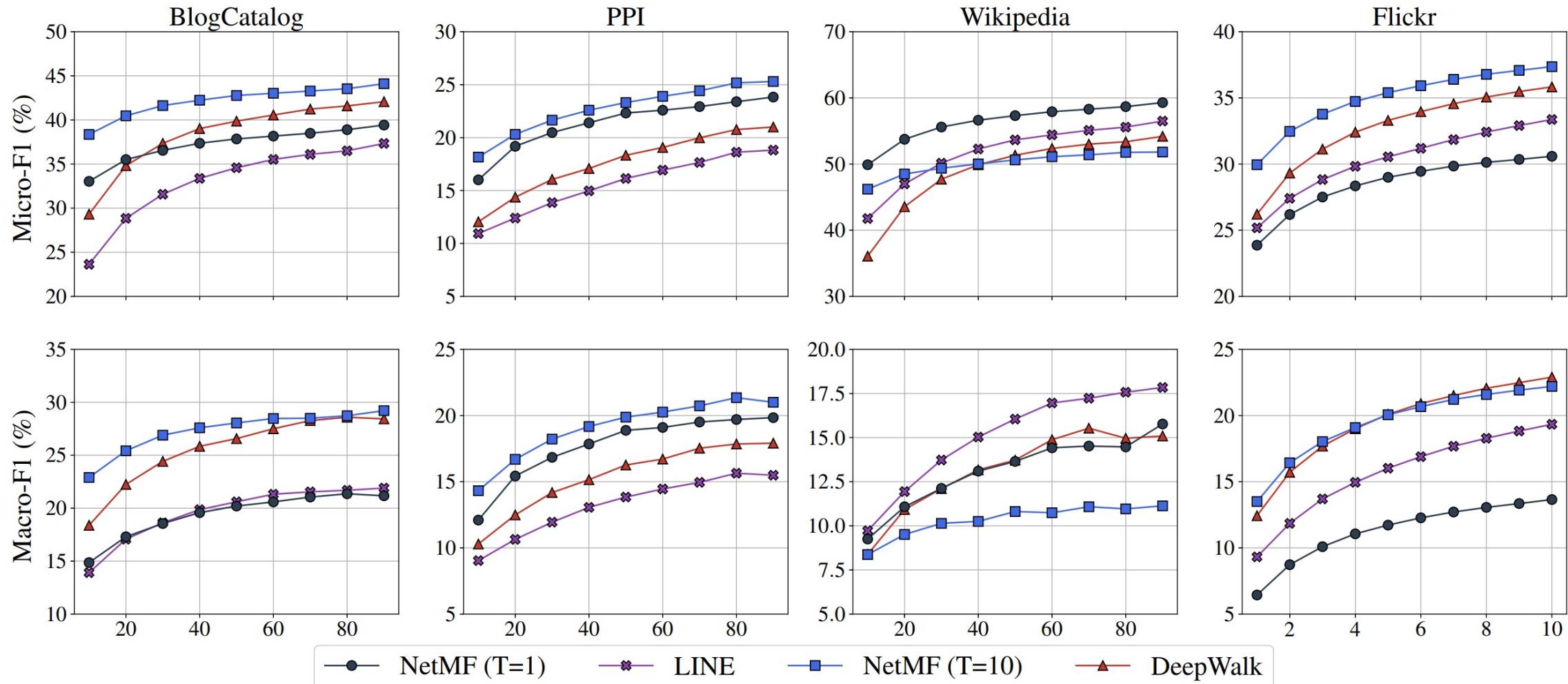
- and because  $M'_{i,j} = \max(M_{i,j}, 1) \geq 1$ , we have

$$\left| M'_{i,j} - \hat{M}'_{i,j} \right| = \left| \max(M_{i,j}, 1) - \max(\hat{M}_{i,j}, 1) \right| \leq \left| M_{i,j} - \hat{M}_{i,j} \right|$$

- Also because the property of NGL,

$$\sigma_s \left( \left( \frac{1}{T} \sum_{r=1}^T P^r \right) D^{-1} \right) \leq \frac{\left| \frac{1}{T} \sum_{r=1}^T \lambda_{p_s}^r \right|}{d_{q_1}} \quad \Rightarrow \quad \| M - \hat{M} \|_F \leq \frac{\text{vol}(G)}{bd_{\min}} \sqrt{\sum_{j=k+1}^n \left| \frac{1}{T} \sum_{r=1}^T \lambda_j^r \right|^2}$$

# Experimental Results

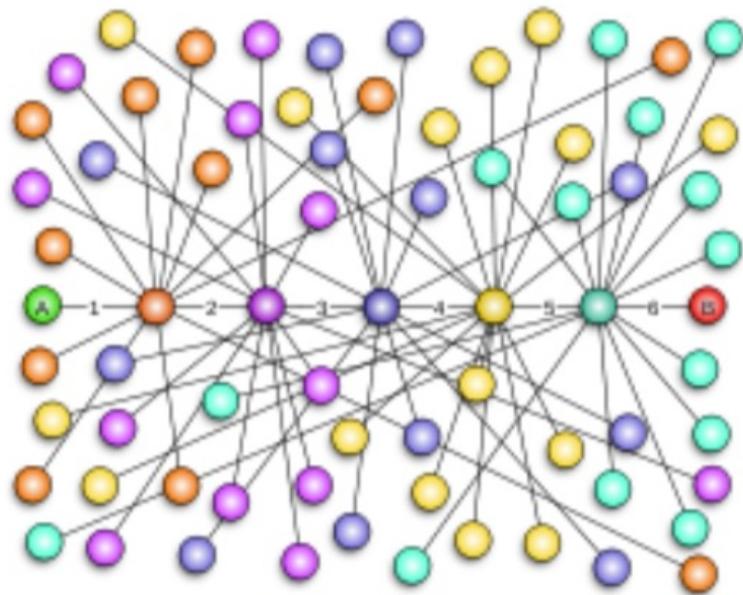


Predictive performance on varying the ratio of training data;  
The x-axis represents the ratio of labeled data (%)

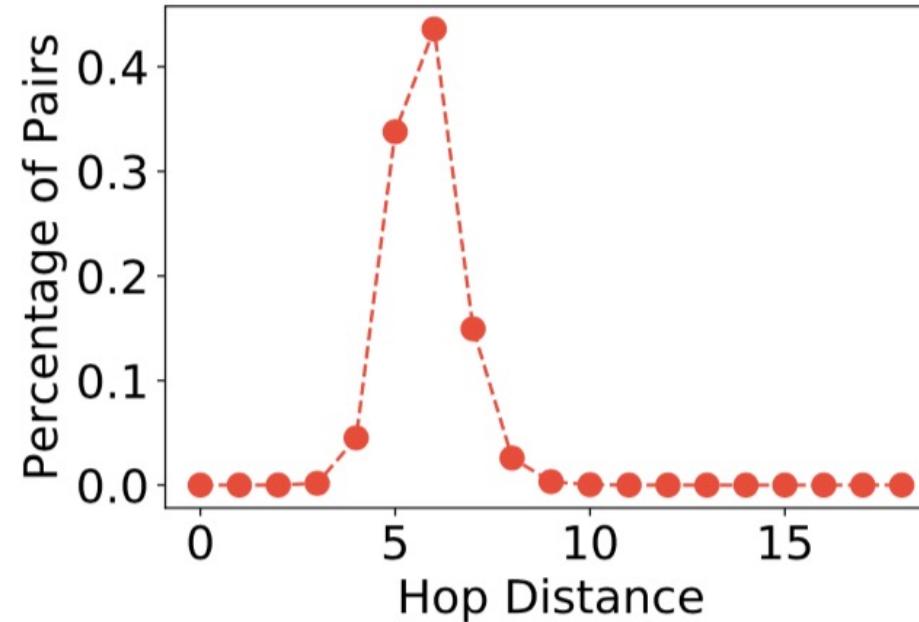
# Network Embedding as Matrix Factorization

- DeepWalk  $\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$
- LINE  $\log \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right)$
- PTE  $\log \begin{pmatrix} \left[ \begin{array}{c} \alpha \text{vol}(G_{\text{ww}}) (\mathbf{D}_{\text{row}}^{\text{ww}})^{-1} \mathbf{A}_{\text{ww}} (\mathbf{D}_{\text{col}}^{\text{ww}})^{-1} \\ \beta \text{vol}(G_{\text{dw}}) (\mathbf{D}_{\text{row}}^{\text{dw}})^{-1} \mathbf{A}_{\text{dw}} (\mathbf{D}_{\text{col}}^{\text{dw}})^{-1} \\ \gamma \text{vol}(G_{\text{lw}}) (\mathbf{D}_{\text{row}}^{\text{lw}})^{-1} \mathbf{A}_{\text{lw}} (\mathbf{D}_{\text{col}}^{\text{lw}})^{-1} \end{array} \right] \end{pmatrix} - \log b$
- node2vec  $\log \left( \frac{\frac{1}{2T} \sum_{r=1}^T \left( \sum_u \mathbf{X}_{w,u} \underline{\mathbf{P}}_{c,w,u}^r + \sum_u \mathbf{X}_{c,u} \underline{\mathbf{P}}_{w,c,u}^r \right)}{b (\sum_u \mathbf{X}_{w,u}) (\sum_u \mathbf{X}_{c,u})} \right)$

# Challenge in NetMF



Small world



Academic graph

$\log^\circ \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$  is always a dense matrix.

# Sparsify $S$

For random-walk matrix polynomial  $L = D - \sum_{r=1}^T \alpha_r D (D^{-1} A)^r$

where  $\sum_{r=1}^T \alpha_r = 1$  and  $\alpha_r$  non-negative

One can construct a  $(1 + \epsilon)$ -spectral sparsifier  $\tilde{L}$  with  $O(n \log n \epsilon^{-2})$  non-zeros

in time  $O(T^2 m \epsilon^{-2} \log^2 n)$

$O(T^2 m \epsilon^{-2} \log n)$  for undirected graphs

1. D. Cheng, Y. Cheng, Y. Liu, R. Peng, and S.H. Teng, Efficient Sampling for Gaussian Graphical Models via Spectral Sparsification, COLT 2015.
2. D. Cheng, Y. Cheng, Y. Liu, R. Peng, and S.H. Teng. Spectral sparsification of random-walk matrix polynomials. arXiv:1502.03496.

# Sparsify $S$

For random-walk matrix polynomial  $L = D - \sum_{r=1}^T \alpha_r D (D^{-1} A)^r$

where  $\sum_{r=1}^T \alpha_r = 1$  and  $\alpha_r$  non-negative

One can construct a **( $1 + \epsilon$ )-spectral sparsifier  $\tilde{L}$**  with  $O(n \log n \epsilon^{-2})$  non-zeros

in time  $O(T^2 m \epsilon^{-2} \log^2 n)$

Suppose  $G = (V, E, A)$  and  $\tilde{G} = (V, \tilde{E}, \tilde{A})$  are two weighted undirected networks. Let  $L = D_G - A$  and  $\tilde{L} = D_{\tilde{G}} - \tilde{A}$  be their Laplacian matrices, respectively. We define  $G$  and  $\tilde{G}$  are  $(1 + \epsilon)$ -spectrally similar if

$$\forall x \in \mathbb{R}^n, (1 - \epsilon) \cdot x^\top \tilde{L} x \leq x^\top L x \leq (1 + \epsilon) \cdot x^\top \tilde{L} x.$$

1. D. Cheng, Y. Cheng, Y. Liu, R. Peng, and S.H. Teng, Efficient Sampling for Gaussian Graphical Models via Spectral Sparsification, COLT 2015.
2. D. Cheng, Y. Cheng, Y. Liu, R. Peng, and S.H. Teng. Spectral sparsification of random-walk matrix polynomials. arXiv:1502.03496.

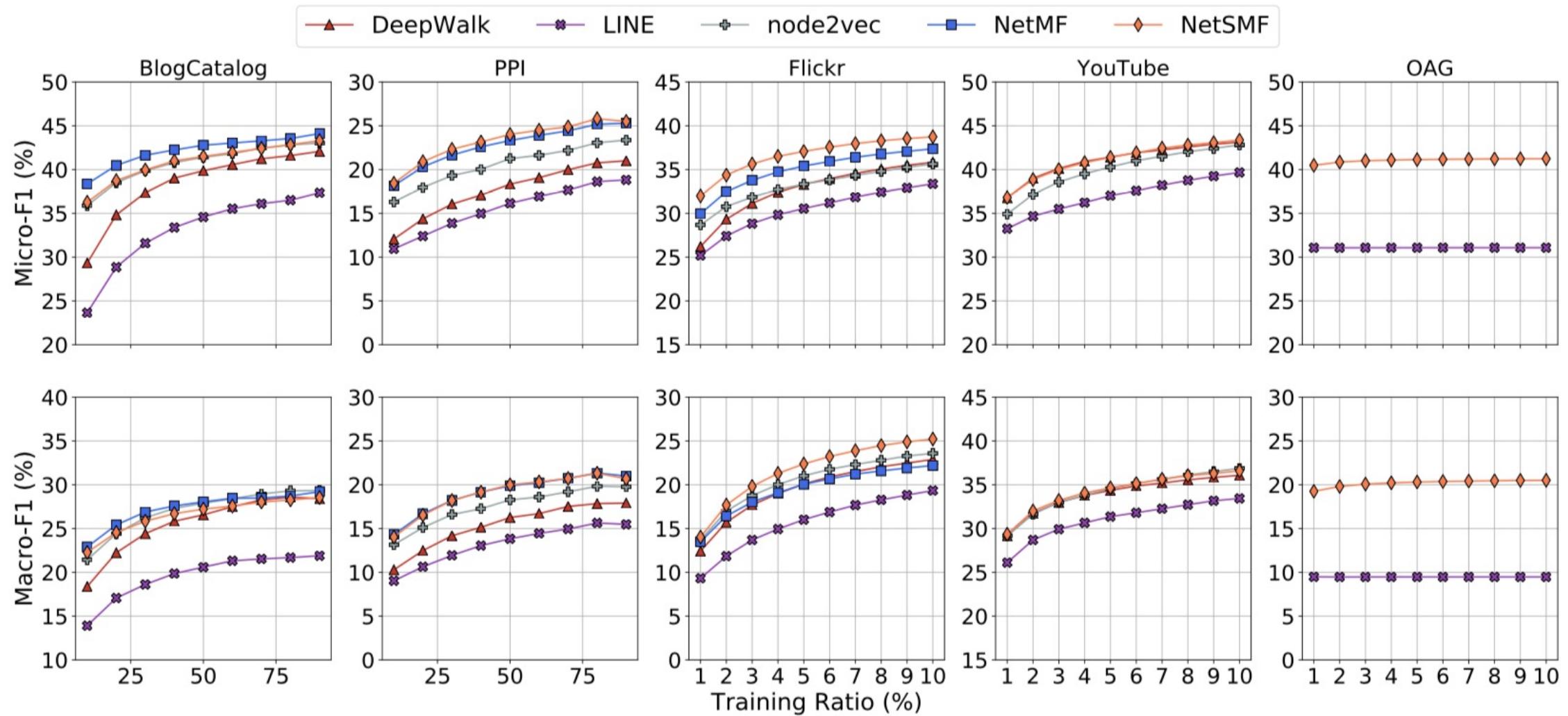
# NetSMF --- Sparse

- ▶ Construct a random walk matrix polynomial sparsifier,  $\tilde{\mathbf{L}}$
- ▶ Construct a NetMF matrix sparsifier.

$$\text{trunc\_log}^\circ \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \tilde{\mathbf{L}}) \mathbf{D}^{-1} \right)$$

- ▶ Factorize the constructed matrix

# Results



# NE as Sparse Matrix Factorization

- node-context set  $\mathcal{D} = E$  (**sparsity**)
- To avoid the trivial solution  $(r_i = c_j, r_i^T c_j \rightarrow \infty, s.t. \hat{p} \rightarrow 1)$
- Local negative samples drawn from

$$P_{\mathcal{D},j} \propto \sum_{i:(i,j) \in \mathcal{D}} p_{i,j}$$

- Modify the loss (sum over the edge-->**sparse**)

$$l = - \sum_{(i,j) \in \mathcal{D}} [p_{i,j} \ln \sigma(r_i^T c_j) + \lambda P_{\mathcal{D},j} \ln \sigma(-r_i^T c_j)]$$

# NE as Sparse Matrix Factorization

- Let the partial derivative w.r.t.  $r_i^T c_j$  be zero

$$r_i^T c_j = \ln p_{i,j} - \ln(\lambda P_{D,j}), \quad (v_i, v_j) \in \mathcal{D}$$

- Matrix to be factorized (**sparse**)

$$M_{i,j} = \begin{cases} \ln p_{i,j} - \ln(\lambda P_{D,j}) & , (v_i, v_j) \in \mathcal{D} \\ 0 & , (v_i, v_j) \notin \mathcal{D} \end{cases}$$

# NE as Sparse Matrix Factorization

- Compared with matrix factorization method (e.g., NetMF)

$$\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right) \text{ v.s. } M_{i,j} = \begin{cases} \ln p_{i,j} - \ln(\lambda P_{D,j}) & , (v_i, v_j) \in \mathcal{D} \\ 0 & , (v_i, v_j) \notin \mathcal{D} \end{cases}$$

- Sparsity** (local structure and local negative samples) → much **faster and scalable** (e.g., randomized tSVD,  $O(|E|)$ )
- The optimization (single thread) is much **faster** than SGD used in DeepWalk, LINE, etc. and is still **scalable!!!**
- Challenge: may **lose** high order information!
- Improvement via **spectral propagation**

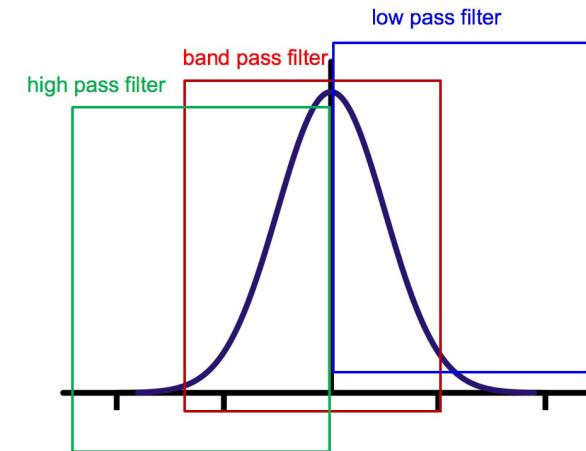
# NE Enhancement via Spectral Propagation

- the form of the spectral filter

$$\tilde{L} = U \text{diag}([g(\lambda_1), \dots, g(\lambda_n)]) U^T$$

$$g(\lambda) = e^{-\frac{1}{2}[(\lambda - \mu)^2 - 1]\theta}$$

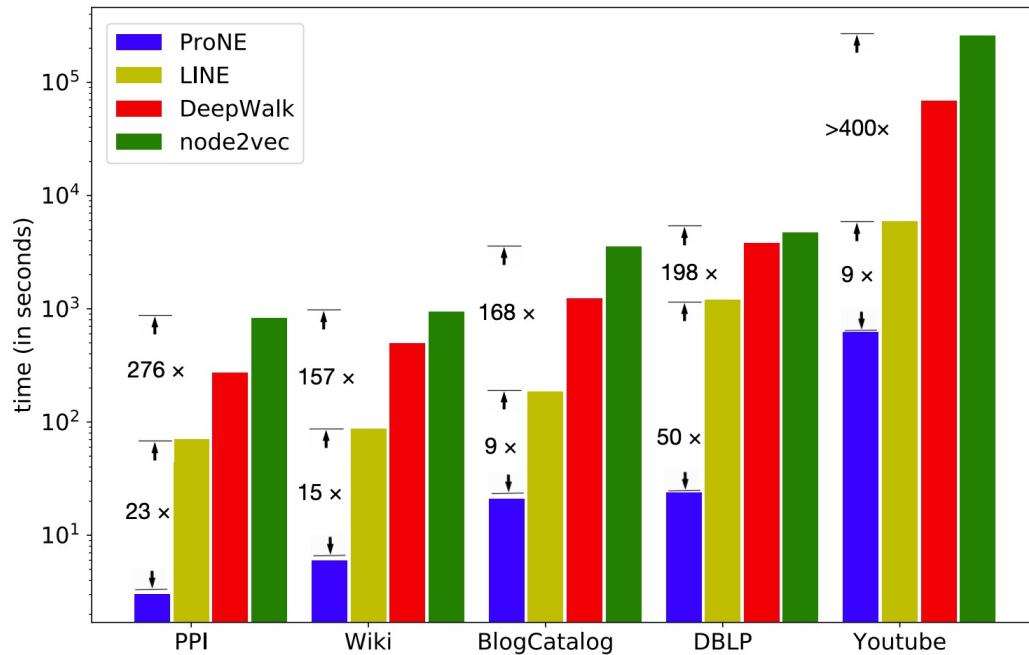
- Band-pass (low-pass, high-pass)
- pass eigenvalues within a certain range and weaken eigenvalues outside that range
- amplify **local** and **global** network information



# Complexity of ProNE

- Spectral propagation only involves sparse matrix multiplication! The complexity is linear!
- Sparse matrix factorization + spectral propagation =  $O(|V|d^2 + k|E|)$

# Results



\* ProNE (1 thread) v.s.  
Others (20 threads)

\* 10 minutes on  
Youtube (~1M nodes)

Dataset	DeepWalk	LINE	node2vec	ProNE
PPI	272	70	828	3
Wiki	494	87	939	6
BlogCatalog	1,231	185	3,533	21
DBLP	3,825	1,204	4,749	24
YouTube	68,272	5,890	>5days	627

# Effectiveness experiments

Dataset	training ratio	0.1	0.3	0.5	0.7	0.9
PPI	DeepWalk	16.4	19.4	21.1	22.3	22.7
	LINE	16.3	20.1	21.5	22.7	23.1
	node2vec	16.2	19.7	21.6	23.1	24.1
	GraRep	15.4	18.9	20.2	20.4	20.9
	HOPE	16.4	19.8	21.0	21.7	22.5
	ProNE (SMF)	15.8	20.6	22.7	23.7	24.2
	ProNE	<b>18.2</b>	<b>22.7</b>	<b>24.6</b>	<b>25.4</b>	<b>25.9</b>
	( $\pm\sigma$ )	( $\pm 0.5$ )	( $\pm 0.3$ )	( $\pm 0.7$ )	( $\pm 1.0$ )	( $\pm 1.1$ )
	DeepWalk	40.4	45.9	48.5	49.1	49.4
	LINE	<b>47.8</b>	50.4	51.2	51.6	52.4
Wiki	node2vec	45.6	47.0	48.2	49.6	50.0
	GraRep	47.2	49.7	50.6	50.9	51.8
	HOPE	38.5	39.8	40.1	40.1	40.1
	ProNE (SMF)	47.6	51.6	53.2	53.5	53.9
	ProNE	47.3	<b>53.1</b>	<b>54.7</b>	<b>55.2</b>	<b>57.2</b>
	( $\pm\sigma$ )	( $\pm 0.7$ )	( $\pm 0.4$ )	( $\pm 0.8$ )	( $\pm 0.8$ )	( $\pm 1.3$ )
	DeepWalk	36.2	39.6	40.9	41.4	42.2
	LINE	28.2	30.6	33.2	35.5	36.8
	BlogCatalog					

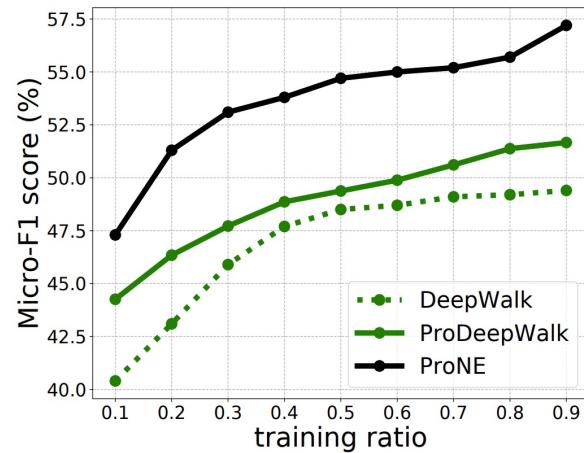
Dataset	training ratio	0.01	0.03	0.05	0.07	0.09
DBLP	DeepWalk	49.3	55.0	57.1	57.9	58.4
	LINE	48.7	52.6	53.5	54.1	54.5
	node2vec	48.9	55.1	57.0	58.0	58.4
	GraRep	50.5	52.6	53.2	53.5	53.8
	HOPE	<b>52.2</b>	55.0	55.9	56.3	56.6
	ProNE (SMF)	50.8	54.9	56.1	56.7	57.0
	ProNE	48.8	<b>56.2</b>	<b>58.0</b>	<b>58.8</b>	<b>59.2</b>
	( $\pm\sigma$ )	( $\pm 1.0$ )	( $\pm 0.5$ )	( $\pm 0.2$ )	( $\pm 0.2$ )	( $\pm 0.1$ )
	DeepWalk	38.0	40.1	41.3	42.1	42.8
	LINE	33.2	35.5	37.0	38.2	39.3
Youtube	ProNE (SMF)	36.5	40.2	41.2	41.7	42.1
	ProNE	<b>38.2</b>	<b>41.4</b>	<b>42.3</b>	<b>42.9</b>	<b>43.3</b>
	( $\pm\sigma$ )	( $\pm 0.8$ )	( $\pm 0.3$ )	( $\pm 0.2$ )	( $\pm 0.2$ )	( $\pm 0.2$ )

\* ProNE (SMF) = ProNE w/

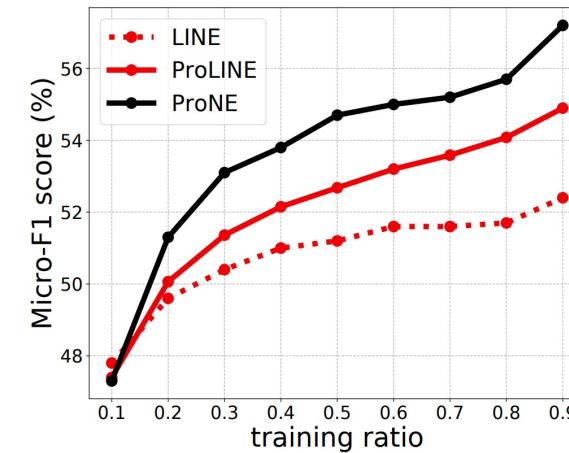
Embed 100,000,000 nodes by one thread:  
29 hours with performance superiority

n

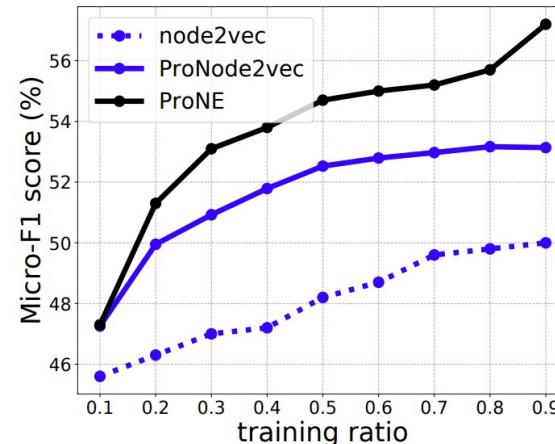
# Spectral Propagation: $k$ -way Cheeger



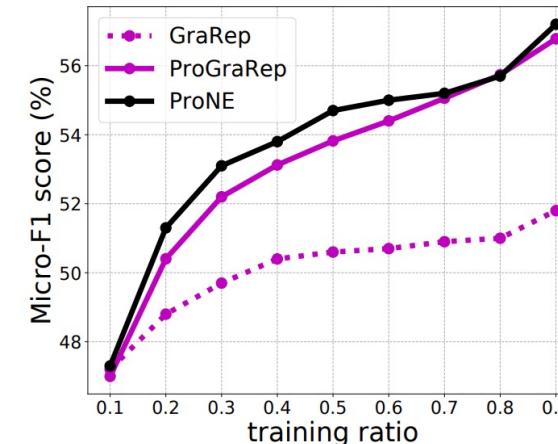
(a) ProDeepWalk



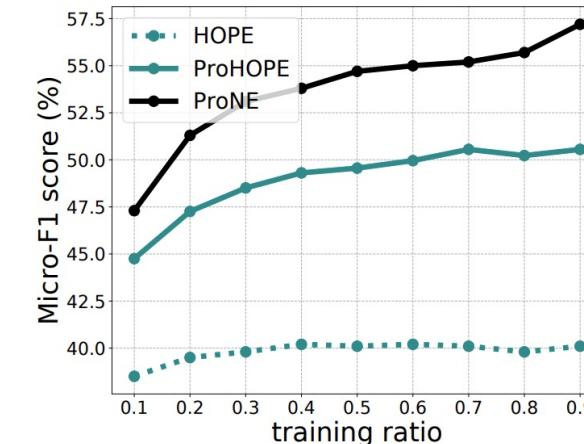
(b) ProLINE



(c) ProNode2vec



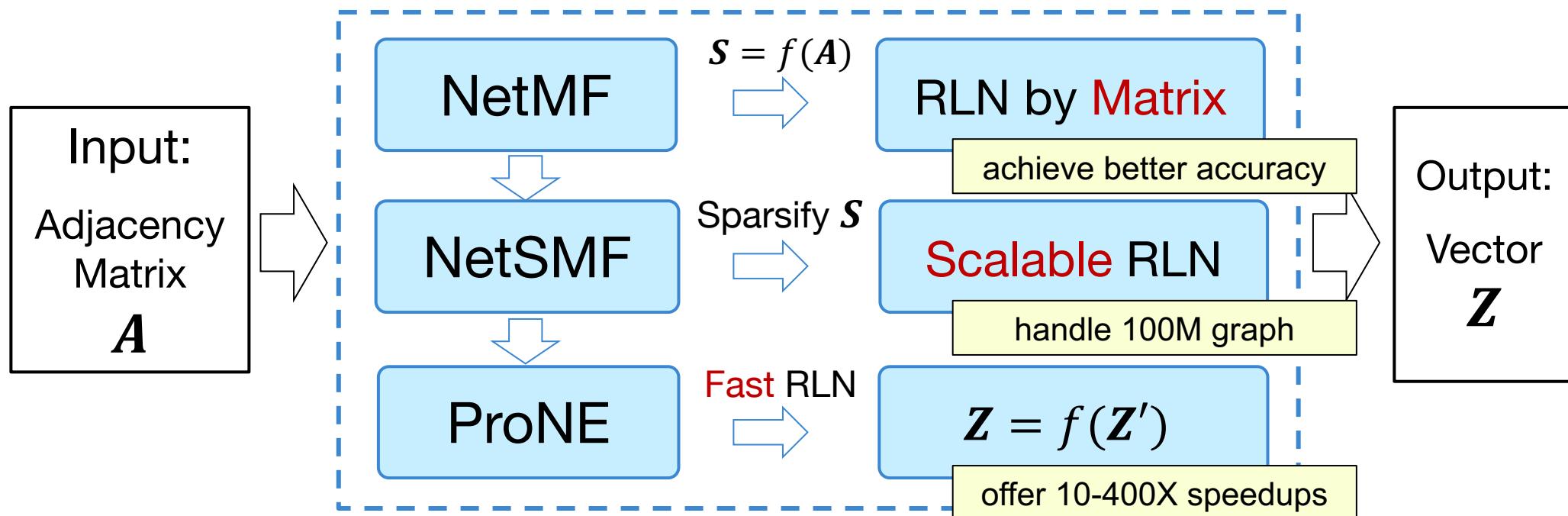
(d) ProGraRep



(e) ProHOPE

$$\frac{\lambda_k}{2} \leq \rho_G(k) \leq O(k^2) \sqrt{\lambda_k}$$

# Representation Learning on Networks



1. Qiu et al. Network embedding as matrix factorization: unifying deepwalk, line, pte, and node2vec. *WSDM'18*. **The most cited paper in WSDM'18 as of May 2019**
2. J. Qiu, Y. Dong, H. Ma, J. Li, C. Wang, K. Wang, and J. Tang. NetSMF: Large-Scale Network Embedding as Sparse Matrix Factorization. *WWW'19*.
3. J. Zhang, Y. Dong, Y. Wang, J. Tang, and M. Ding. ProNE: Fast and Scalable Network Representation Learning. *IJCAI'19*.

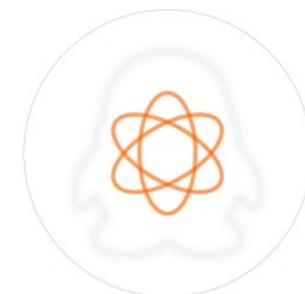
# GRL: NE&GNN

Network  
Embedding

Matrix  
Factorization

Graph Neural  
Networks

GNN  
Pre-Training



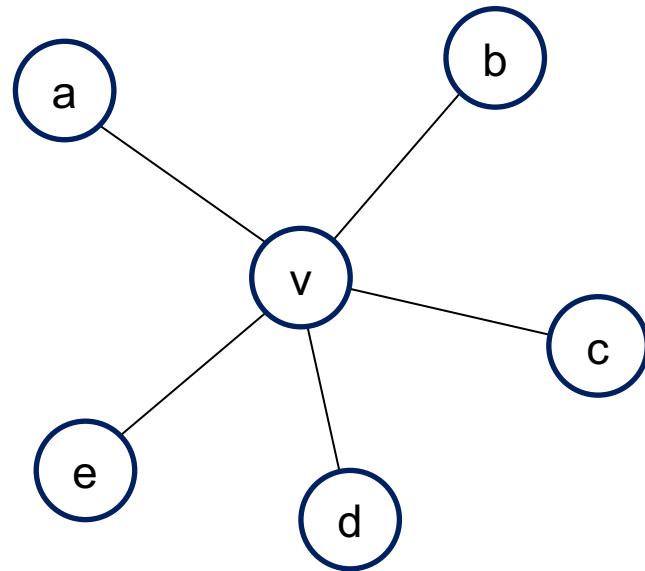
<https://alchemy.tencent.com/>



CogDL

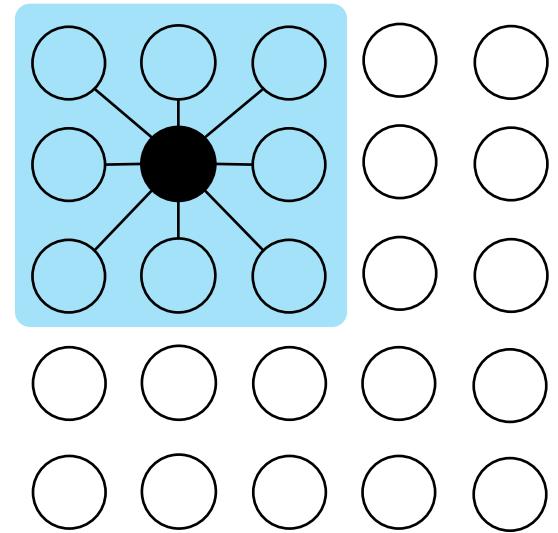
<https://github.com/thudm/cogdl>

# Graph Neural Networks



**Graph Convolution**

1. Choose neighborhood
2. Determine the order of selected neighbors
3. Parameter sharing



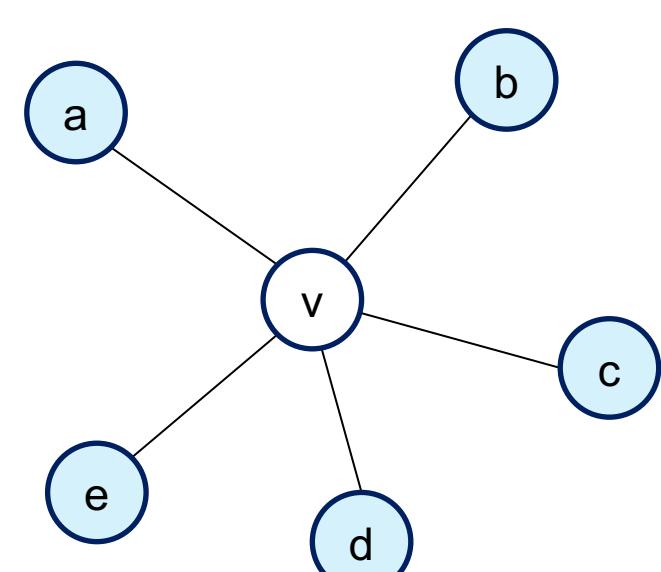
**CNN**

## Neighborhood Aggregation:

- Aggregate neighbor information and pass into a neural network
- It can be viewed as a center-surround filter in CNN---graph convolutions!

1. Niepert et al. Learning Convolutional Neural Networks for Graphs. In ICML 2016
2. Defferrard et al. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In NIPS 2016
3. Huawei Shen. Graph Neural Networks. A video talk (in Chinese), June 2020. [https://dl.ccf.org.cn/audioVideo/detail.html?id=4966039790962688&\\_ack=1](https://dl.ccf.org.cn/audioVideo/detail.html?id=4966039790962688&_ack=1)

# Graph Convolutional Networks



Non-linear activation function (e.g., ReLU)

parameters in layer  $k$

$$h_v^k = \sigma(W^k)$$

node  $v$ 's embedding at layer  $k$

$$\sum_{u \in N(v) \cup v} \frac{h_u^{k-1}}{\sqrt{|N(u)||N(v)|}})$$

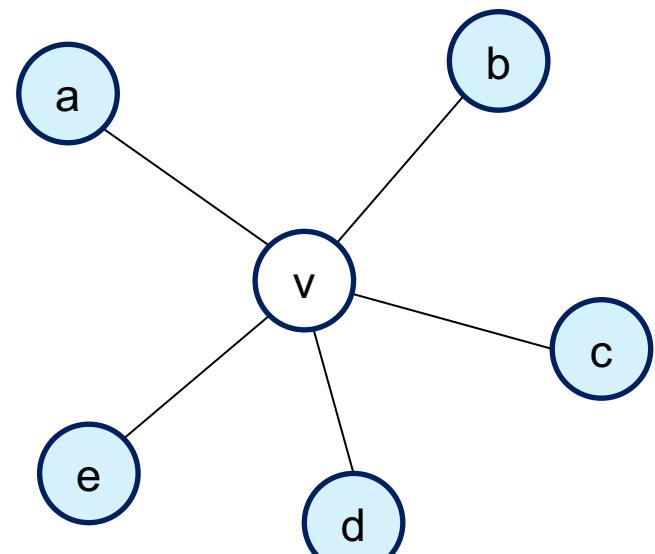
the neighbors of node  $v$

$$H^k = \sigma(\hat{A} H^{(k-1)} W^{(k)})$$

normalized Laplacian matrix

Aggregate info from neighborhood via the normalized Laplacian matrix

# Graph Convolutional Networks

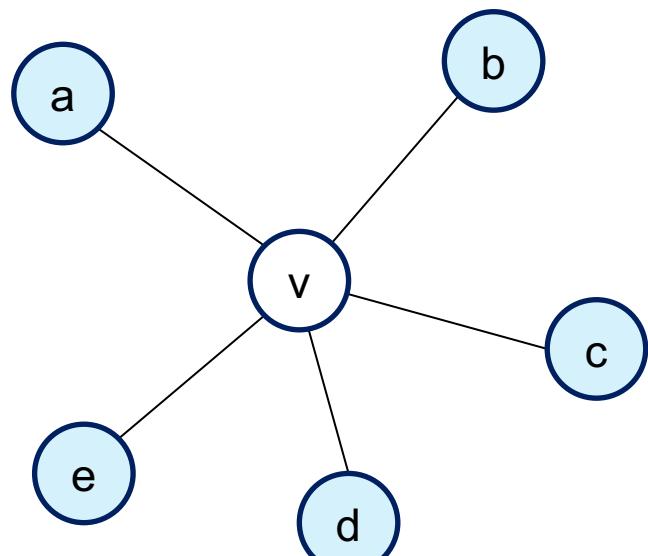


$$h_v^k = \sigma(W^k \sum_{u \in N(v)} \frac{h_u^{k-1}}{\sqrt{|N(u)||N(v)|}} + W^k \sum_v \frac{h_v^{k-1}}{\sqrt{|N(v)||N(v)|}})$$

Aggregate from itself

Aggregate from  $v$ 's neighbors

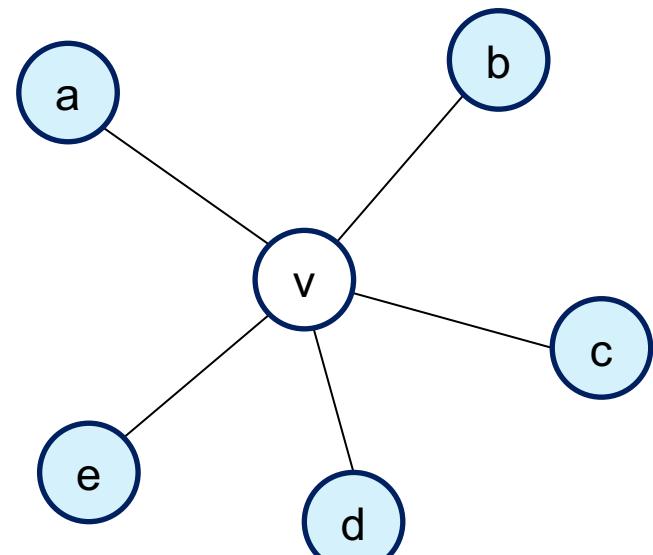
# Graph Convolutional Networks



The same parameters for both its neighbors & itself

$$h_v^k = \sigma(W^k \sum_{u \in N(v)} \frac{h_u^{k-1}}{\sqrt{|N(u)||N(v)|}} + W^k \sum_v \frac{h_v^{k-1}}{\sqrt{|N(v)||N(v)|}})$$

# Graph Convolutional Networks

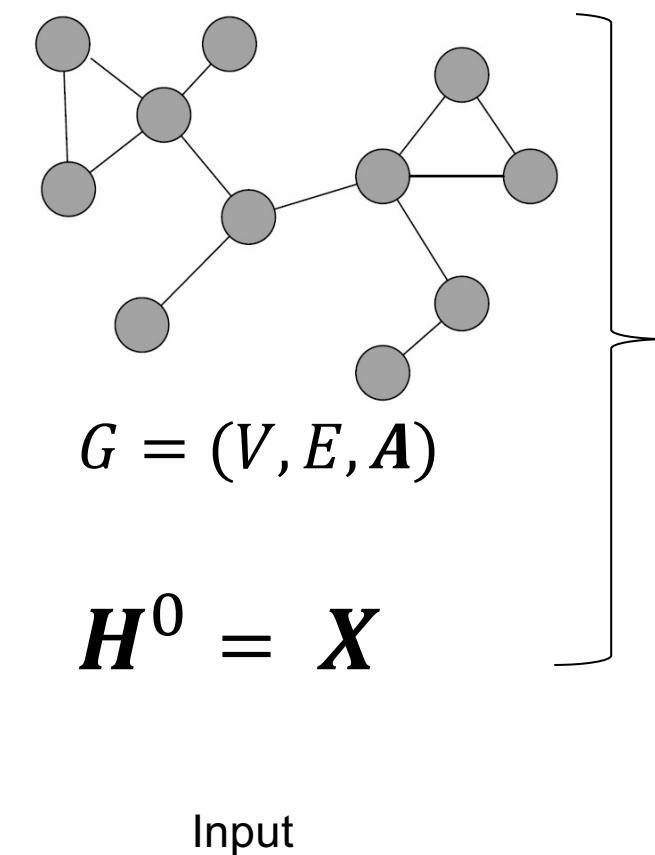


$$\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)}$$

$$\mathbf{h}_v^k = \sigma(\mathbf{W}^k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}} + \mathbf{W}^k \sum_v \frac{\mathbf{h}_v^{k-1}}{\sqrt{|N(v)||N(v)|}})$$

$$\mathbf{D}^{-\frac{1}{2}} \mathbf{I} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)}$$

# Graph Convolutional Networks

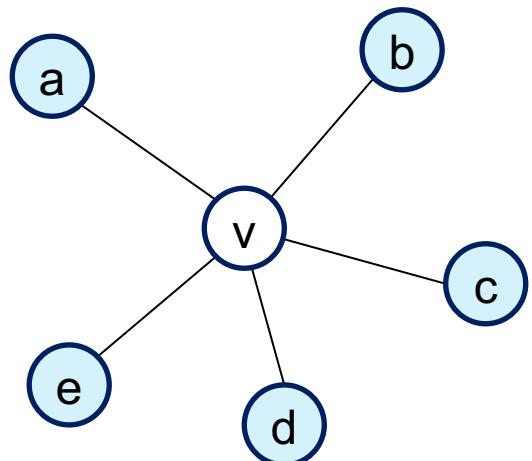


$$H^k = \sigma \left( D^{-\frac{1}{2}} (A + I) D^{-\frac{1}{2}} H^{(k-1)} W^{(k)} \right)$$

$$\Rightarrow Z = H^K$$

- Model training
  - The common setting is to have an end to end training framework with a supervised task
- Benefits: Parameter sharing for all nodes
  - #parameters is sublinear in |V|
  - Enable inductive learning for new nodes

# GraphSAGE



GCN

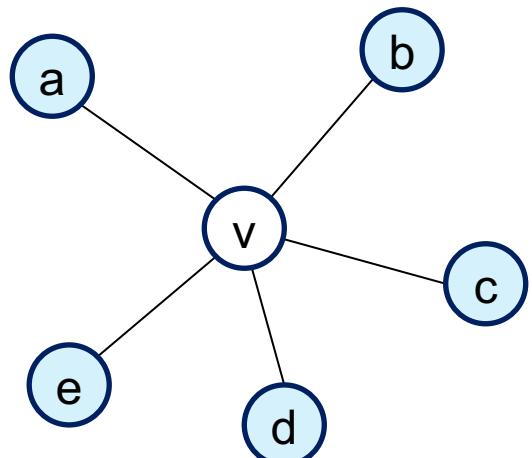
$$\mathbf{h}_v^k = \sigma(\mathbf{W}^k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}})$$

GraphSage

$$\mathbf{h}_v^k = \sigma([A^k \cdot \text{AGG}(\{\mathbf{h}_u^{k-1}, \forall u \in N(v)\}), \mathbf{B}^k \mathbf{h}_v^{k-1}])$$

**Generalized aggregation:** any differentiable function that maps set of vectors to a single vector

# GraphSAGE



GCN

$$\mathbf{h}_v^k = \sigma(\mathbf{W}^k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}})$$

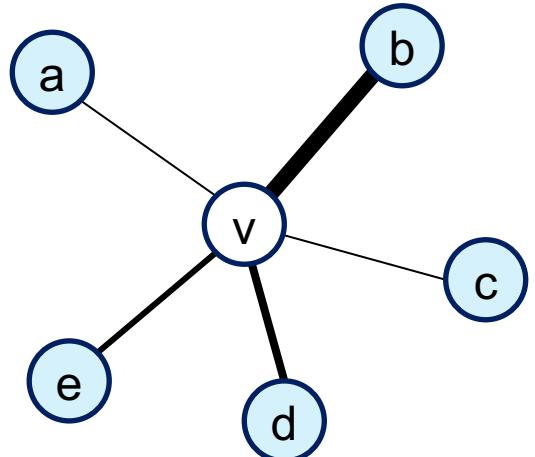
GraphSage

Instead of summation, it concatenates  
neighbor & self embeddings

$$\mathbf{h}_v^k = \sigma([\mathbf{A}^k \cdot \text{AGG}(\{\mathbf{h}_u^{k-1}, \forall u \in N(v)\}), \mathbf{B}^k \mathbf{h}_v^{k-1}])$$

**Generalized aggregation:** any differentiable  
function that maps set of vectors to a single vector

# GNN: Graph Attention



GCN

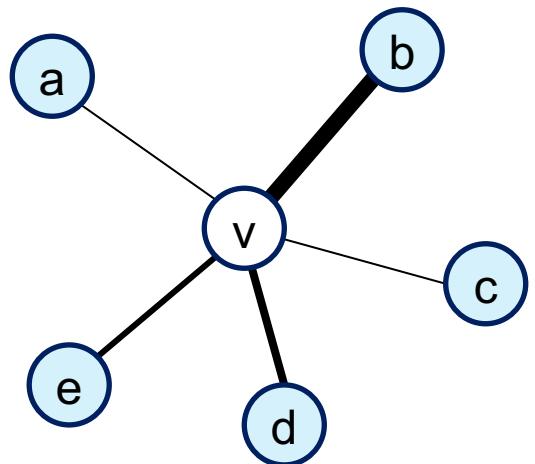
$$\mathbf{h}_v^k = \sigma(\mathbf{W}^k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}})$$

Graph Attention

$$\mathbf{h}_v^k = \sigma\left( \sum_{u \in N(v) \cup v} \alpha_{v,u} \mathbf{W}^k \mathbf{h}_u^{k-1} \right)$$

Learned attention weights

# GNN: Graph Attention



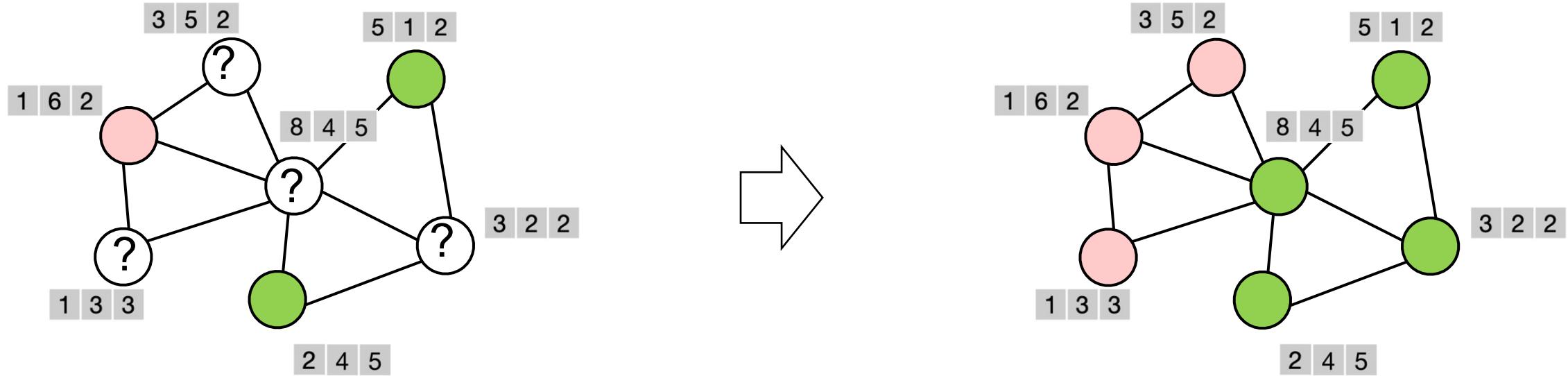
$$\alpha_{v,u} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{Q}\mathbf{h}_v, \mathbf{Q}\mathbf{h}_u]))}{\sum_{u' \in N(v) \cup \{v\}} \exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{Q}\mathbf{h}_v, \mathbf{Q}\mathbf{h}_{u'}]))}$$

Various ways to define attention!

# Revisiting Semi-Supervised Learning on Graphs

- W. Feng, J. Zhang, Y. Dong, Y. Han, H. Luan, Q. Xu, Q. Yang, E. Kharlamov, and J. Tang. Graph Random Neural Networks for Semi-Supervised Learning on Graphs. NeurIPS'20.  
<https://arxiv.org/abs/2005.11079>
- Code & data for Grand: <https://github.com/Grand20/grand>

# Semi-Supervised Learning on Graphs



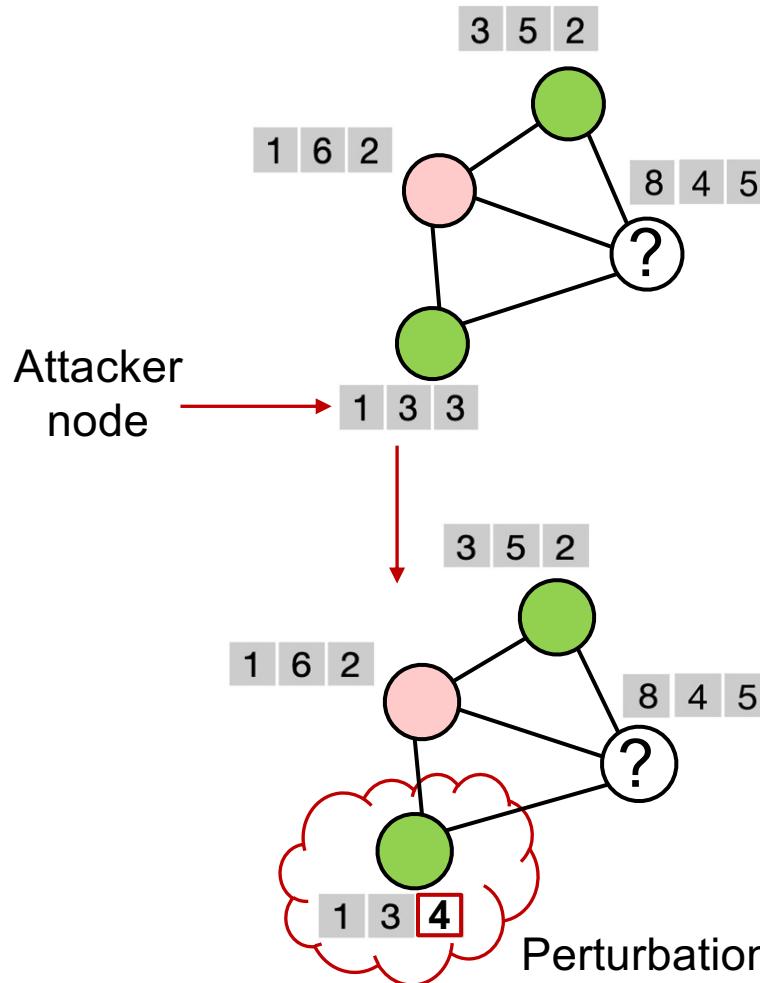
**Input:** a partially labeled & attributed graph

**Output:** infer the labels of unlabeled nodes

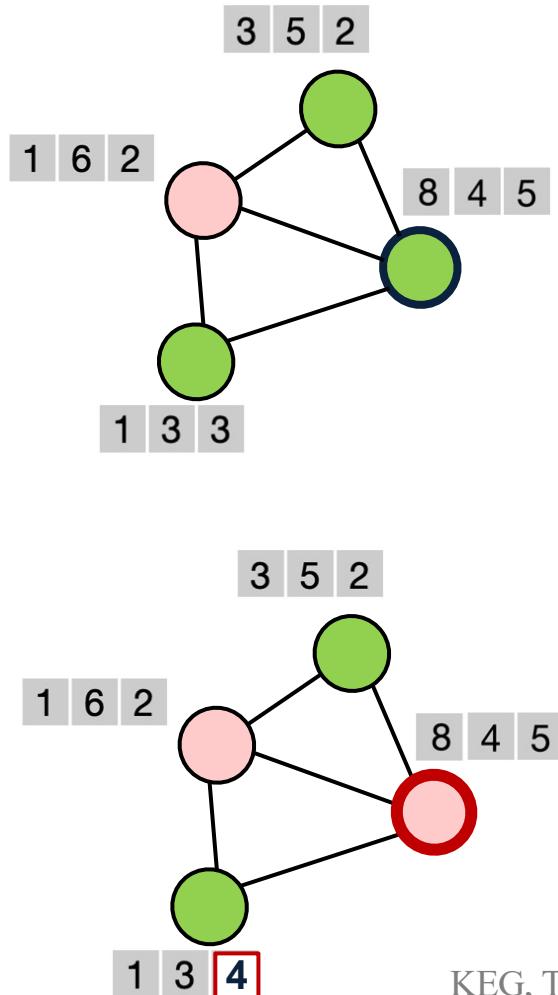
# Graph Neural Networks

$$\mathbf{H}^{k+1} = \sigma(\hat{\mathbf{A}}\mathbf{H}^{(k)}\mathbf{W}^{(k)})$$

a deterministic propagation



1. Each node is highly dependent with its neighbors, making GNNs **non-robust** to noises



# Graph Neural Networks

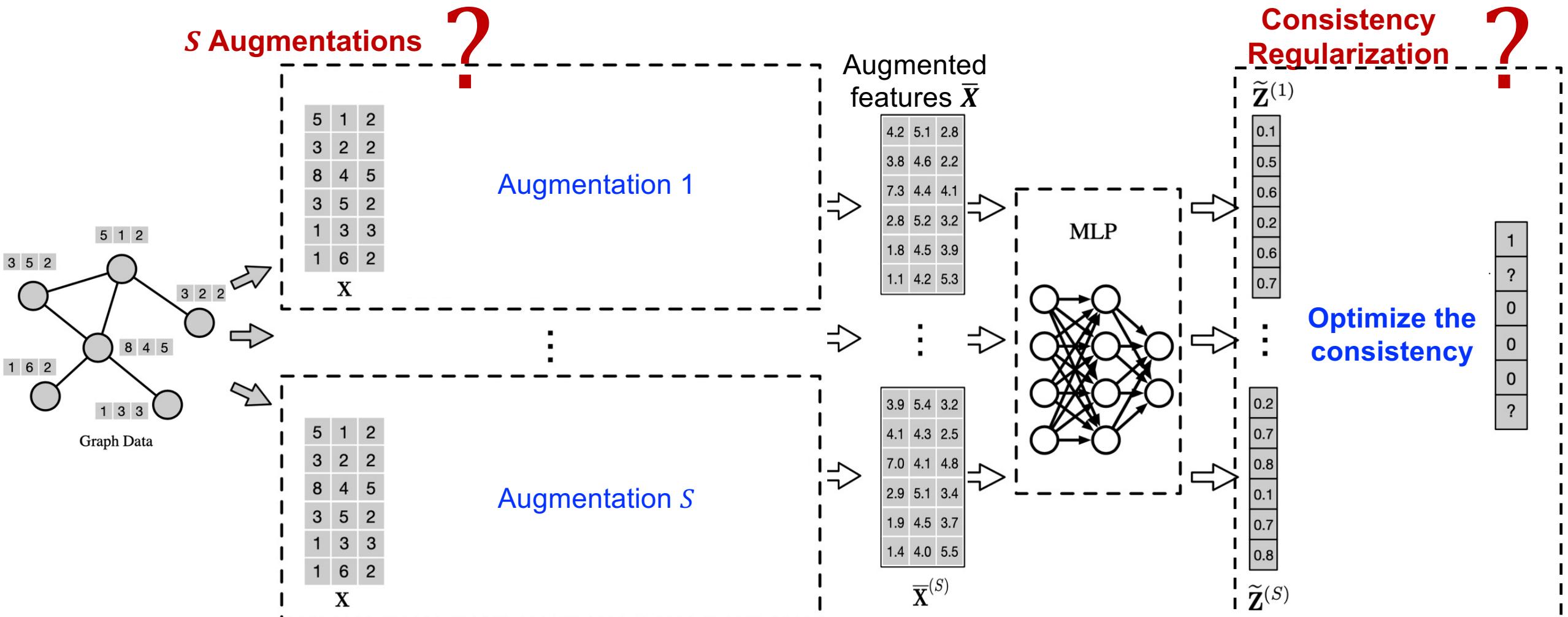
$$H^{k+1} = \sigma(\widehat{A}H^{(k)}W^{(k)})$$

feature propagation  
is coupled with  
non-linear transformation

1. Each node is highly dependent with its neighbors, making GNNs **non-robust** to noises
2. Stacking many GNN layers may cause **over-fitting** & **over-smoothing**.

# Graph Random Neural Networks (GRAND)

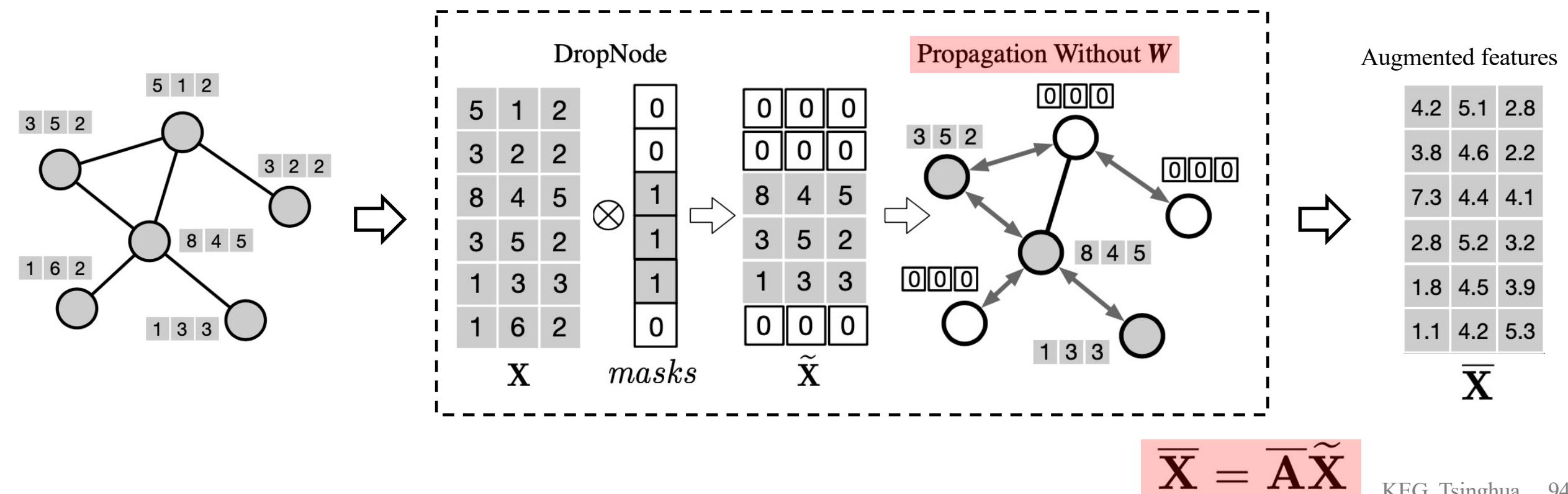
- Consistency Regularized Training:
  - Generates  $S$  data augmentations of the graph
  - Optimizing the consistency among  $S$  augmentations of the graph.



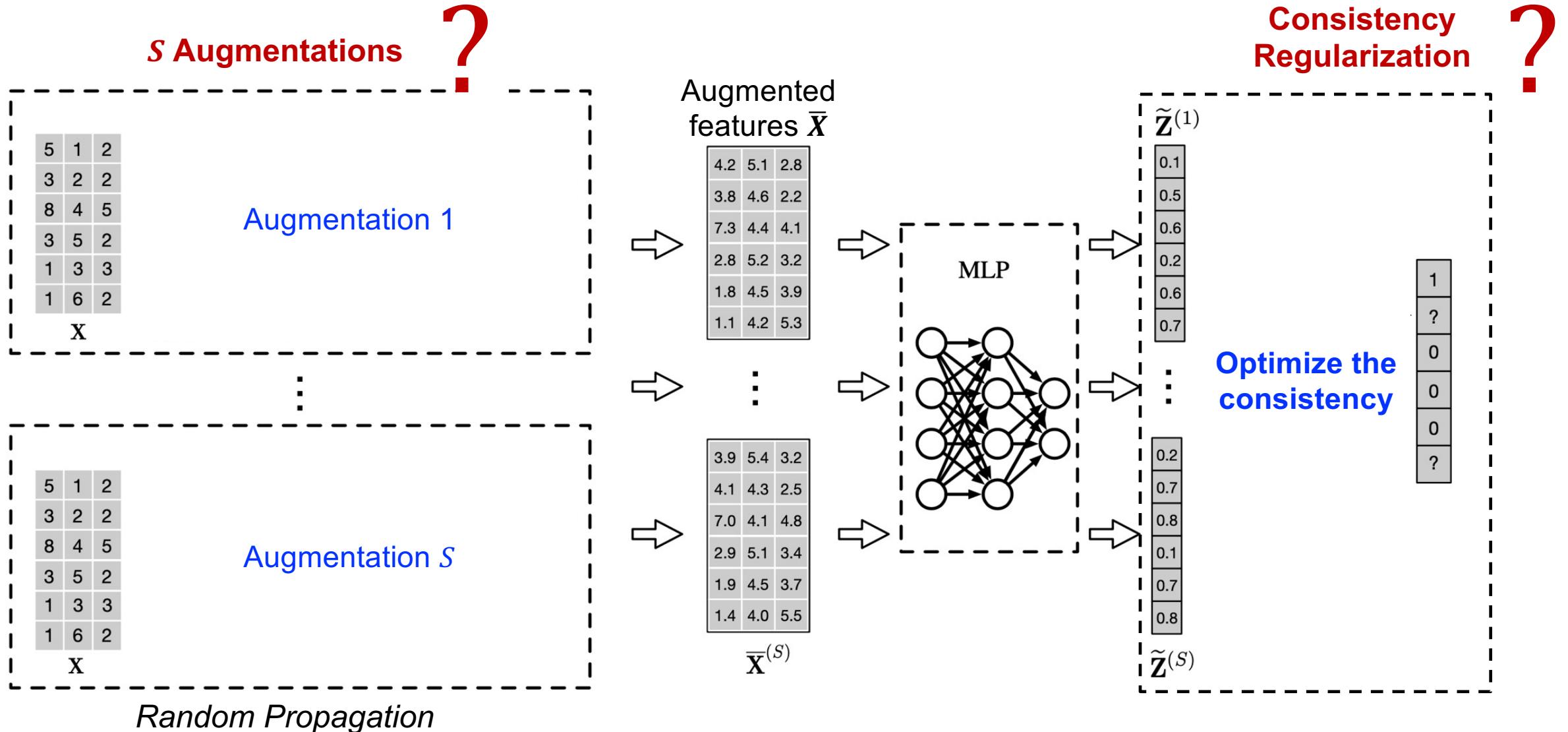
# GRAND: Random Propagation for Graph Data Augmentation

- **Random Propagation** (DropNode + Propagation):
  - Each node is enabled to be not sensitive to specific neighborhoods.
  - Decouple feature propagation from feature transformation.

*Random Propagation for Augmentation*

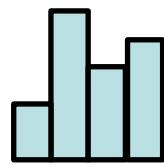


# Graph Random Neural Networks (GRAND)

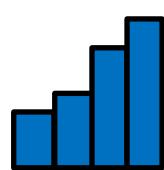
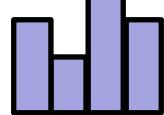


# GRAND: Consistency Regularization

Distributions of a node  
after augmentations



Average



$$\bar{\mathbf{Z}}_i = \frac{1}{S} \sum_{s=1}^S \tilde{\mathbf{Z}}_i^{(s)}$$

$$\mathcal{L}_{sup} = -\frac{1}{S} \sum_{s=1}^S \sum_{i=0}^{m-1} \mathbf{Y}_i^\top \log \tilde{\mathbf{Z}}_i^{(s)}$$

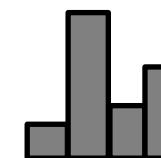


$$\mathcal{L} = \mathcal{L}_{sup} + \lambda \mathcal{L}_{con}$$

$$\mathcal{L}_{con} = \frac{1}{S} \sum_{s=1}^S \sum_{i=0}^{n-1} \mathcal{D}(\bar{\mathbf{Z}}_i', \tilde{\mathbf{Z}}_i^{(s)})$$

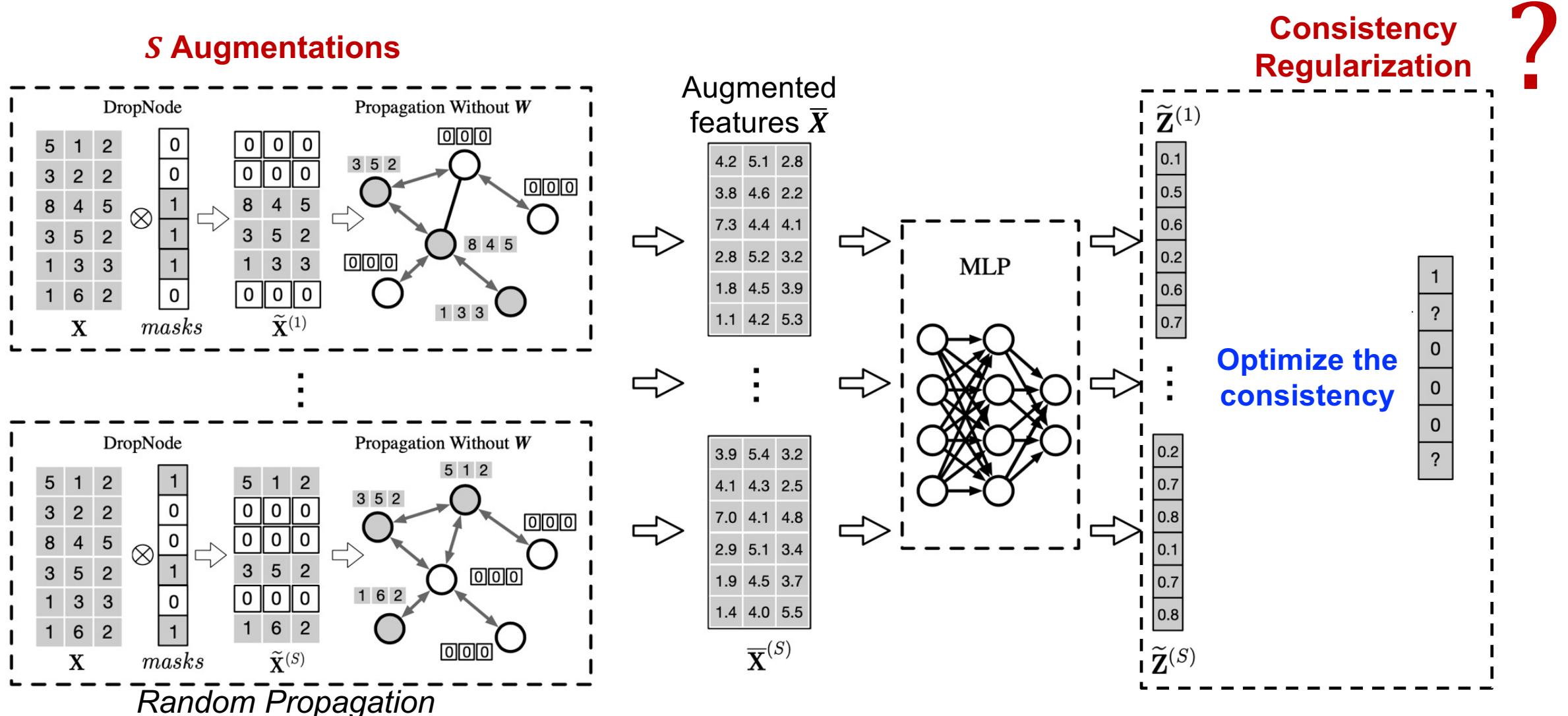


Sharpening



$$\bar{\mathbf{Z}}_{ik}' = \bar{\mathbf{Z}}_{ik}^{\frac{1}{T}} \Bigg/ \sum_{j=0}^{C-1} \bar{\mathbf{Z}}_{ij}^{\frac{1}{T}}$$

# Graph Random Neural Networks (GRAND)



# Consistency Regularized Training Algorithm

## Input:

Adjacency matrix  $\hat{\mathbf{A}}$ , feature matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , times of augmentations in each epoch  $S$ , DropNode probability  $\delta$ .

## Output:

Prediction  $\mathbf{Z}$ .

- 1: **while** not convergence **do**
- 2:   **for**  $s = 1 : S$  **do**
- 3:     Apply DropNode via Algorithm 1:  $\tilde{\mathbf{X}}^{(s)} \sim \text{DropNode}(\mathbf{X}, \delta)$ .
- 4:     Perform propagation:  $\bar{\mathbf{X}}^{(s)} = \frac{1}{K+1} \sum_{k=0}^K \hat{\mathbf{A}}^k \tilde{\mathbf{X}}^{(s)}$ .
- 5:     Predict class distribution using MLP:  $\tilde{\mathbf{Z}}^{(s)} = P(\mathbf{Y} | \bar{\mathbf{X}}^{(s)}; \Theta)$ .
- 6:   **end for**
- 7:   Compute supervised classification loss  $\mathcal{L}_{sup}$  via Eq. 4 and consistency regularization loss via Eq. 6.
- 8:   Update the parameters  $\Theta$  by gradients descending:  
$$\nabla_{\Theta} \mathcal{L}_{sup} + \lambda \mathcal{L}_{con}$$

Generate  
 $S$  Augmentations

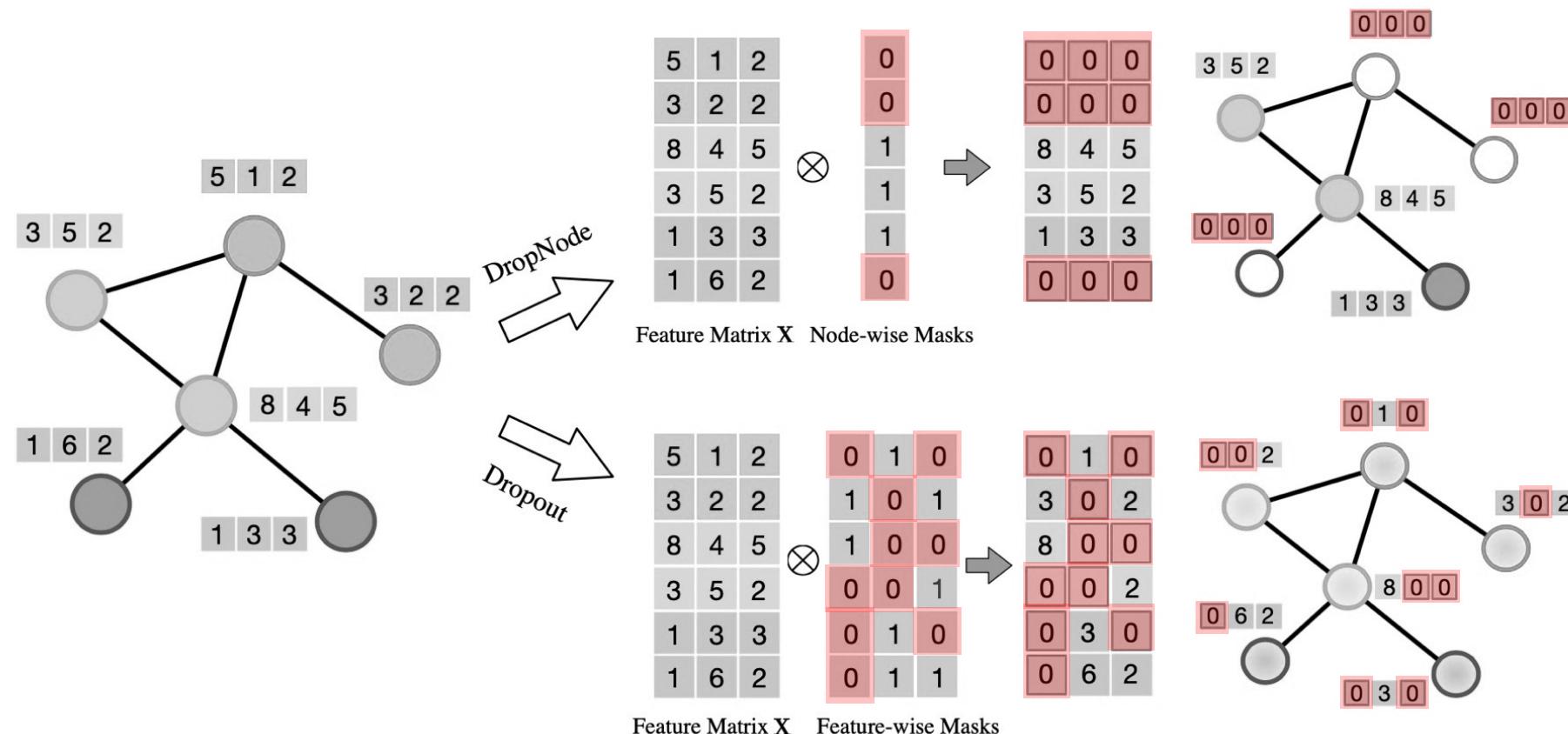
Consistency  
Regularization

- 9: **end while**

- 10: Output prediction  $\mathbf{Z}$  via Eq. 8.

# GRAND: DropNode vs Dropout

- Dropout drops each element in  $X$  independently
- DropNode drops the entire features of selected nodes, i.e., the row vectors of  $X$ , randomly
  - Theoretically, Dropout is an adaptive  $L_2$  regularization.



# Graph Random Neural Networks (GRAND)

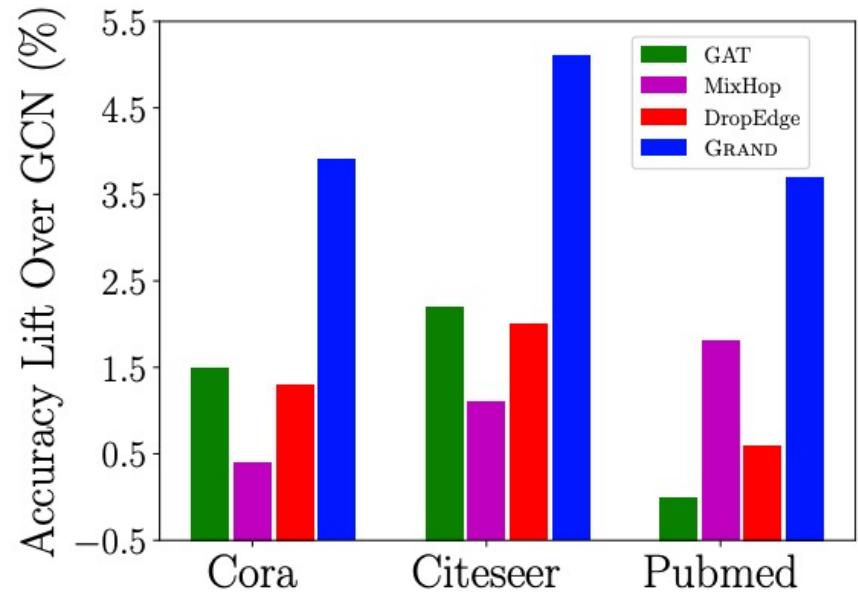
- With Consistency Regularization Loss:
  - Random propagation can enforce the consistency of the classification confidence between each node and its all multi-hop neighborhoods.

$$\begin{aligned}\mathbb{E}_\epsilon(\mathcal{L}_{con}) \approx \mathcal{R}^c(\mathbf{W}) &= \sum_{i=0}^{n-1} z_i^2(1-z_i)^2 \text{Var}_\epsilon \left( \overline{\mathbf{A}}_i \tilde{\mathbf{X}} \cdot \mathbf{W} \right) \\ \mathcal{R}_{DN}^c(\mathbf{W}) &= \frac{\delta}{1-\delta} \sum_{j=0}^{n-1} \left[ (\mathbf{X}_j \cdot \mathbf{W})^2 \sum_{i=0}^{n-1} (\overline{\mathbf{A}}_{ij})^2 z_i^2 (1-z_i)^2 \right] \\ \mathcal{R}_{Do}^c(\mathbf{W}) &= \frac{\delta}{1-\delta} \sum_{h=0}^{d-1} \mathbf{W}_h^2 \sum_{j=0}^{n-1} \left[ \mathbf{X}_{jh}^2 \sum_{i=0}^{n-1} z_i^2 (1-z_i)^2 (\overline{\mathbf{A}}_{ij})^2 \right]\end{aligned}$$

- With Supervised Cross-Entropy Loss:
  - Random propagation can enforce the consistency of the classification confidence between each node and its labeled multi-hop neighborhoods.

# Results

	Method	Cora	Citeseer	Pubmed
GCNs	GCN [19]	81.5	70.3	79.0
	GAT [32]	83.0±0.7	72.5±0.7	79.0±0.3
	APPNP [20]	83.8±0.3	71.6±0.5	79.7±0.3
	Graph U-Net [11]	84.4±0.6	73.2±0.5	79.6±0.2
	SGC [36]	81.0±0.0	71.9±0.1	78.9±0.0
	MixHop [1]	81.9±0.4	71.4±0.8	80.8±0.6
	GMNN [28]	83.7	72.9	81.8
	GraphNAS [12]	84.2±1.0	73.1±0.9	79.6±0.4
Sampling GCNs	GraphSAGE [16]	78.9±0.8	67.4±0.7	77.8±0.6
	FastGCN [7]	81.4±0.5	68.8±0.9	77.6±0.5
Reg. GCNs	VBAT [10]	83.6±0.5	74.0±0.6	79.9±0.4
	G <sup>3</sup> NN [24]	82.5±0.2	74.4±0.3	77.9±0.4
	GraphMix [33]	83.9±0.6	74.5±0.6	81.0±0.6
	DropEdge [29]	82.8	72.3	79.6
	<b>GRAND</b>	<b>85.4±0.4</b>	<b>75.4±0.4</b>	<b>82.7±0.6</b>



Instead of the marginal improvements by conventional GNN baselines over GCN,  
**GRAND achieves much more significant performance lift in all three datasets!**

- W. Feng, J. Zhang, Y. Dong, Y. Han, H. Luan, Q. Xu, Q. Yang, E. Kharlamov, and J. Tang. Graph Random Neural Networks for Semi-Supervised Learning on Graphs. NeurIPS'20. <https://arxiv.org/abs/2005.11079>
- Code & data for Grand: <https://github.com/Grand20/grand>

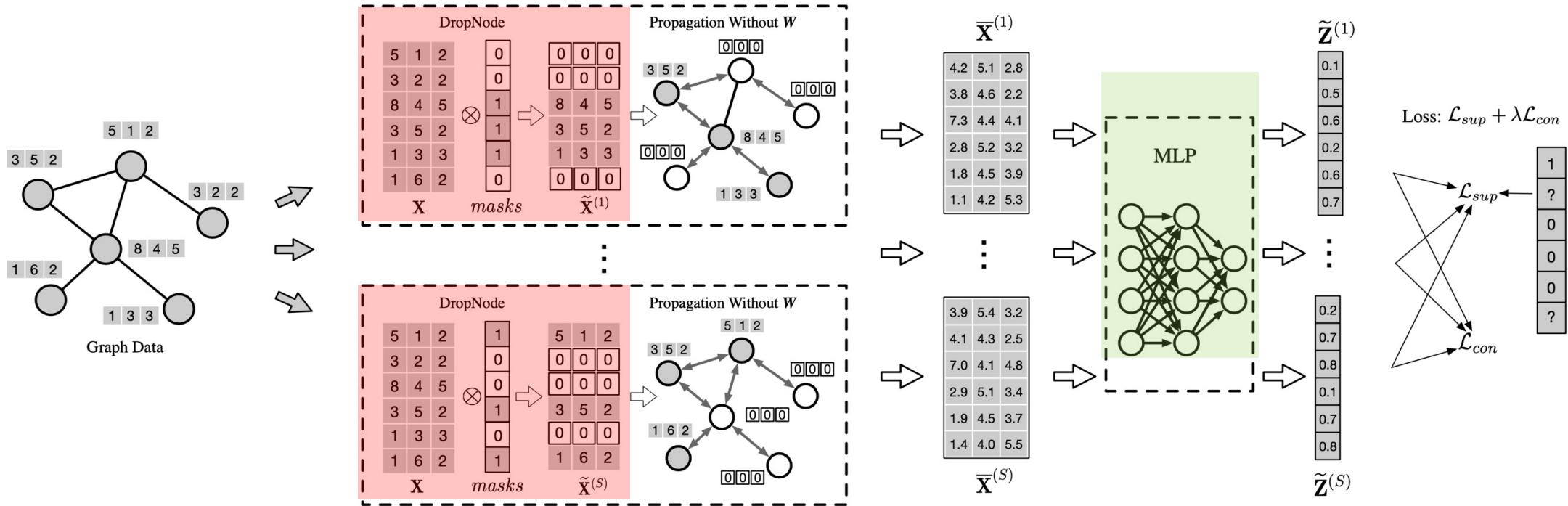
# Results

Table 5: Results on large datasets.

Method	Cora Full	Coauthor CS	Coauthor Physics	Amazon Computer	Amazon Photo	Citation CS
GCN	$62.2 \pm 0.6$	$91.1 \pm 0.5$	$92.8 \pm 1.0$	$82.6 \pm 2.4$	$91.2 \pm 1.2$	$49.9 \pm 2.0$
GAT	$51.9 \pm 1.5$	$90.5 \pm 0.6$	$92.5 \pm 0.9$	$78.0 \pm 19.0$	$85.7 \pm 20.3$	$49.6 \pm 1.7$
GRAND	<b><math>63.5 \pm 0.6</math></b>	<b><math>92.9 \pm 0.5</math></b>	<b><math>94.6 \pm 0.5</math></b>	<b><math>85.7 \pm 1.8</math></b>	<b><math>92.5 \pm 1.7</math></b>	<b><math>52.8 \pm 1.2</math></b>

More experiments on larger graph datasets

# Results



<b>GRAND_dropout</b>	$84.9 \pm 0.4$	$75.0 \pm 0.3$	$81.7 \pm 1.0$
<b>GRAND_GCN</b>	$84.5 \pm 0.3$	$74.2 \pm 0.3$	$80.0 \pm 0.3$
<b>GRAND_GAT</b>	$84.3 \pm 0.4$	$73.2 \pm 0.4$	$79.2 \pm 0.6$
<b>GRAND</b>	<b><math>85.4 \pm 0.4</math></b>	<b><math>75.4 \pm 0.4</math></b>	<b><math>82.7 \pm 0.6</math></b>

Evaluation of the design choices in GRAND

# Results

Method	Cora	Citeseer	Pubmed
GCN [19]	81.5	70.3	79.0
GAT [32]	83.0±0.7	72.5±0.7	79.0±0.3
APPNP [20]	83.8±0.3	71.6±0.5	79.7±0.3
Graph U-Net [11]	84.4±0.6	73.2±0.5	79.6±0.2
SGC [36]	81.0±0.0	71.9±0.1	78.9±0.0
MixHop [1]	81.9±0.4	71.4±0.8	80.8±0.6
GMNN [28]	83.7	72.9	81.8
GraphNAS [12]	84.2±1.0	73.1±0.9	79.6±0.4
DropEdge [29]	82.8	72.3	79.6
w/o CR	84.4±0.5	73.1±0.6	80.9±0.8
w/o mDN	84.7±0.4	74.8±0.4	81.0±1.1
w/o sharpening	84.6±0.4	72.2±0.6	81.6±0.8
w/o CR & DN	83.2±0.5	70.3±0.6	78.5±1.4

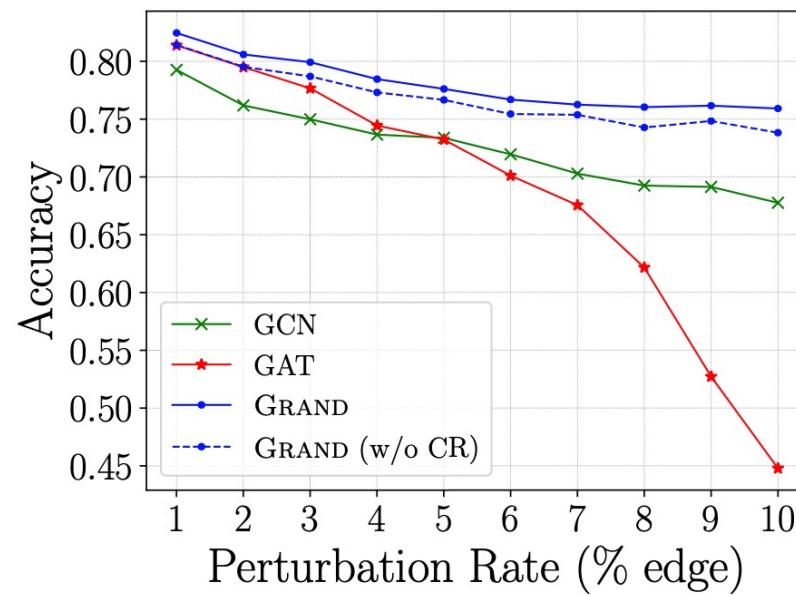
## Ablation Study

1. Each of the designed components contributes to the success of GRAND.
2. GRAND w/o consistency regularization outperforms almost *all 8 non-regularization based GCNs & DropEdge*

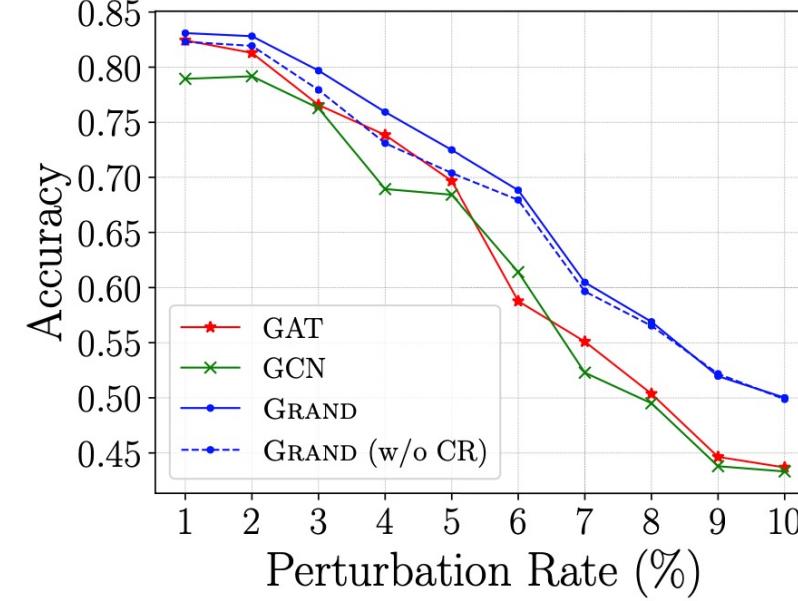
Random Propagation  
vs.

Feature Propagation & Non-Linear Transformation

# Results



(a) Random Attack



(b) Metattack

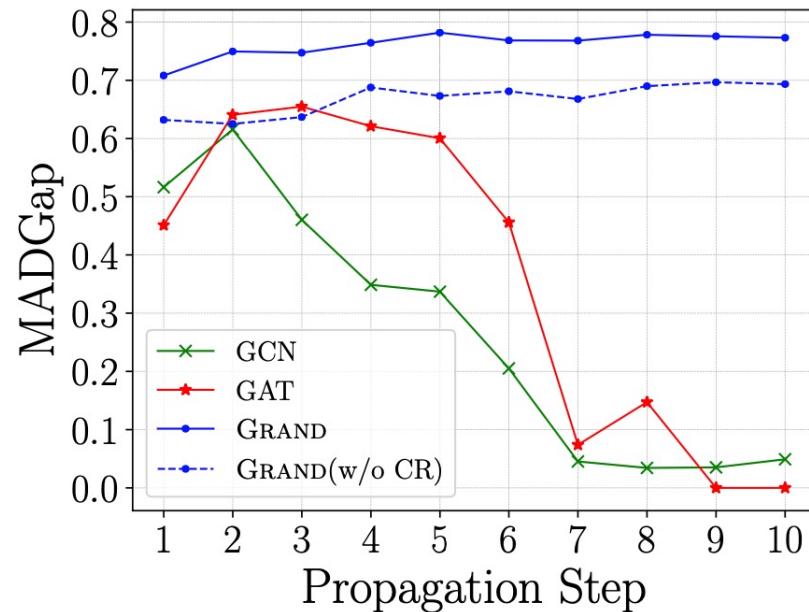
## Robustness

1. GRAND (with or w/o) consistency regularization is more robust than GCN and GAT.

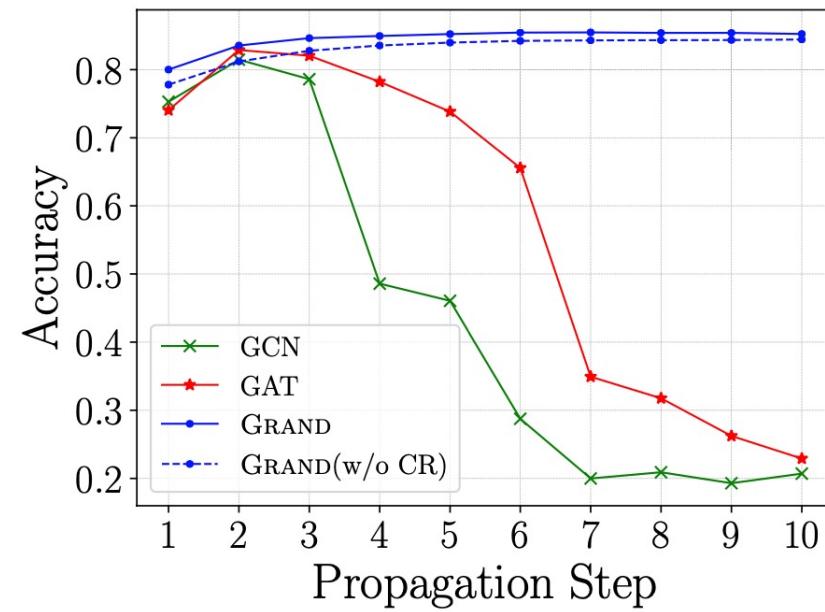
Random Propagation  
vs.

Feature Propagation & Non-Linear Transformation

# Results



(a) MADGap



(b) Classification Results

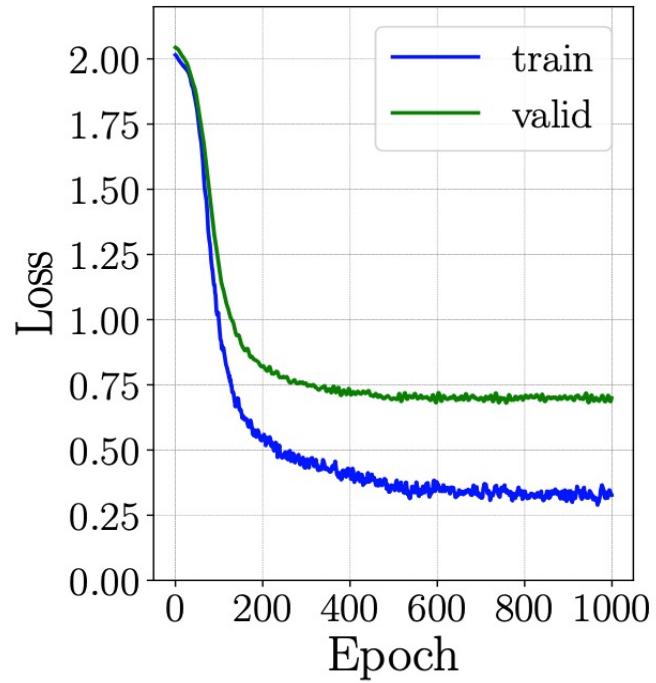
## Over-Smoothing

1. GRAND is powerful to relieve over-smoothing, while GCN & GAT are vulnerable to it

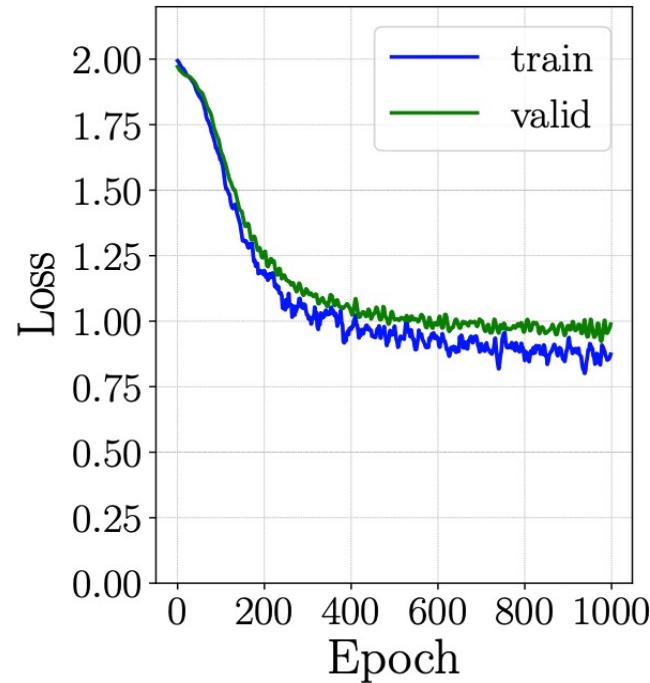
Random Propagation  
vs.

Feature Propagation & Non-Linear Transformation

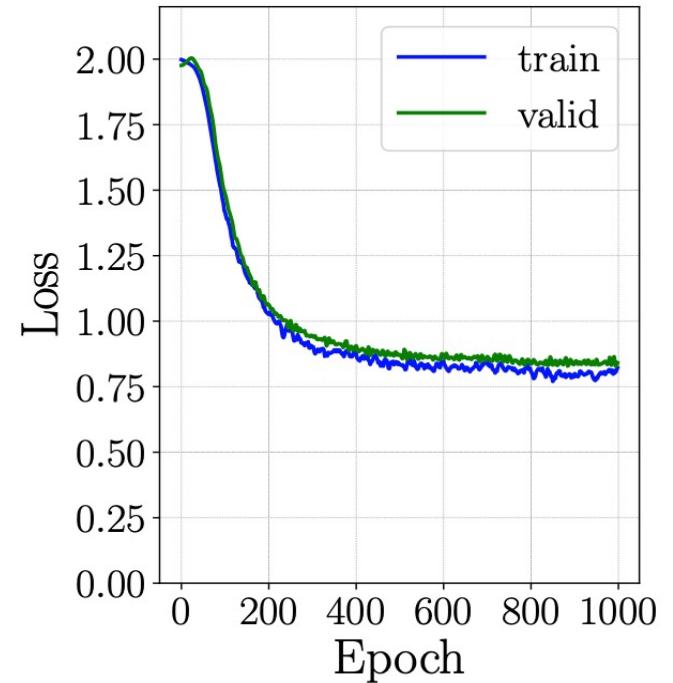
# Results



(a) Without CR and RP



(b) Without CR



(c) GRAND

## Generalization

1. Both the random propagation and consistency regularization improve GRAND's generalization capability

# Graph Robustness Benchmark: Benchmarking the Adversarial Robustness of Graph Machine Learning

**Qinkai Zheng<sup>†</sup>, Xu Zou<sup>†</sup>, Yuxiao Dong<sup>‡\*</sup>, Yukuo Cen<sup>†</sup>,**  
**Da Yin<sup>†</sup>, Jiarong Xu<sup>○</sup>, Yang Yang<sup>◊</sup>, Jie Tang<sup>†‡</sup>**

<sup>†</sup> Department of Computer Science and Technology, Tsinghua University

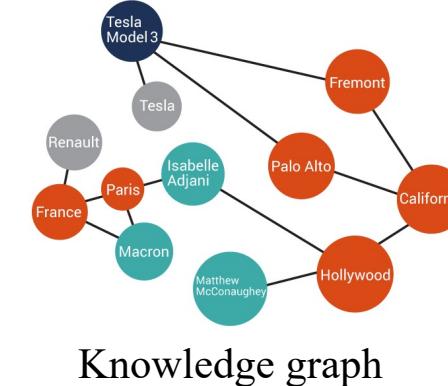
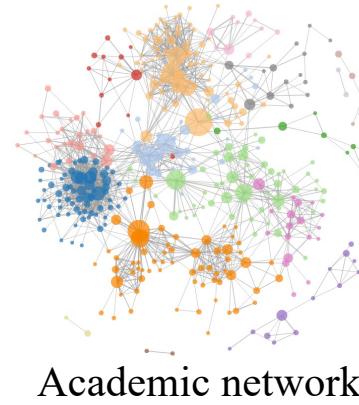
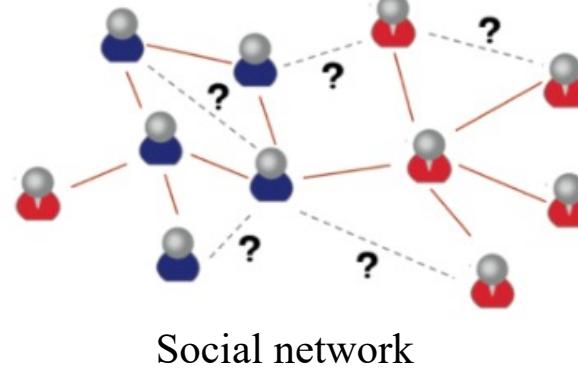
<sup>‡</sup> Microsoft Research, Redmond <sup>○</sup> Fudan University <sup>◊</sup> Zhejiang University



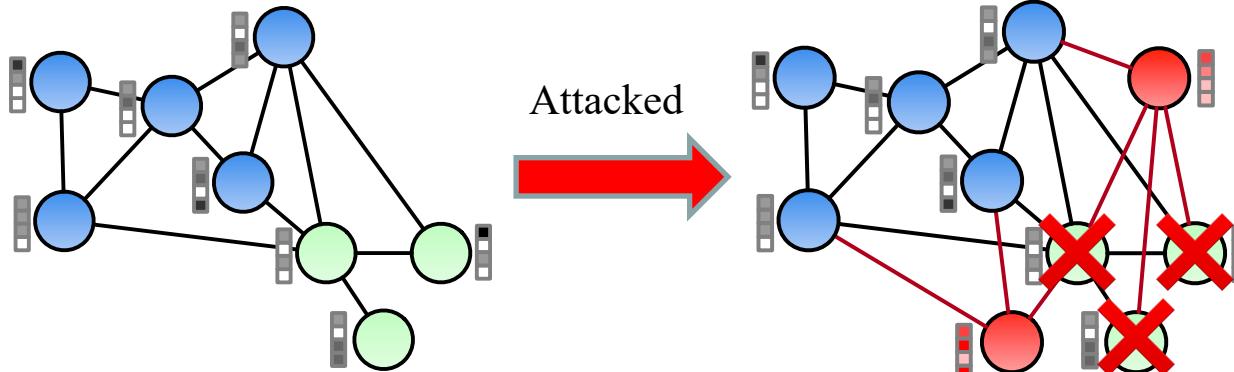
Homepage: <https://cogdl.ai/grb/home> Code: <https://github.com/THUDM/grb>

# Adversarial Robustness in Graph Machine Learning

- Graph machine learning (GML) has made great progress.



- However, GML models face big threats from adversarial attacks.



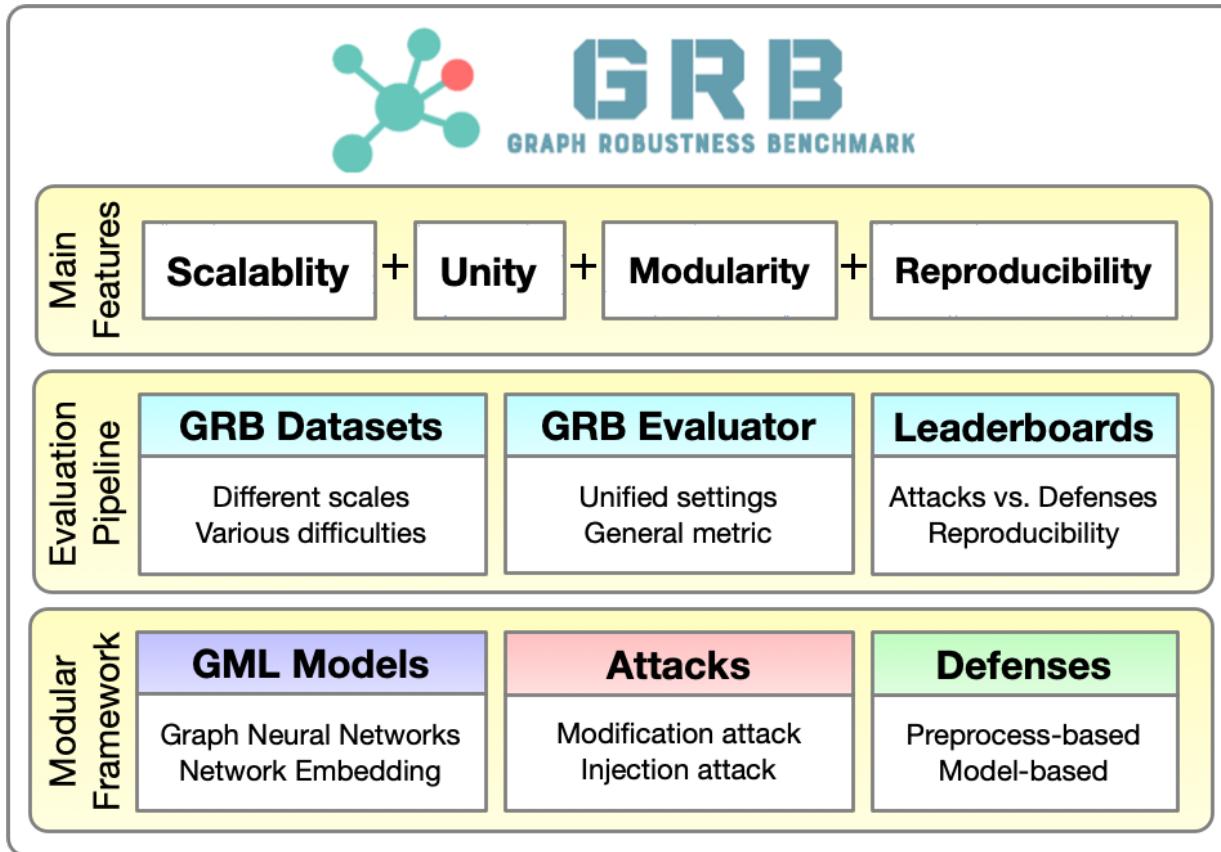
**Adversarial Robustness:**  
*The ability to maintain performance under potential adversarial attacks.*

1. X. Zou, Q. Zheng, Y. Dong, X. Guan, E. Kharlamov, J. Lu, and J. Tang. TDGIA: Effective Injection Attacks on Graph Neural Networks. **KDD'21**.

2. Zheng et al., Graph Robustness Benchmark: Benchmarking the Adversarial Robustness of Graph Machine Learning, **NeurIPS'21**.

Code & data for Grand: <https://cogdl.ai/grb/home> <https://github.com/THUDM/grb>

# Graph Robustness Benchmark (GRB)



- **Elaborated datasets**
  - From small- to large- scale graphs.
  - Robust-specific splitting and normalization.
- **Unified evaluation protocol**
  - Unified settings and general metrics.
  - Unified evaluation pipeline.
- **Refined attack/defense scenarios**
  - More practical scenarios.
  - More realistic assumptions.
- **Modular coding framework**
  - Easy implementation.
  - Assurance of reproducibility.

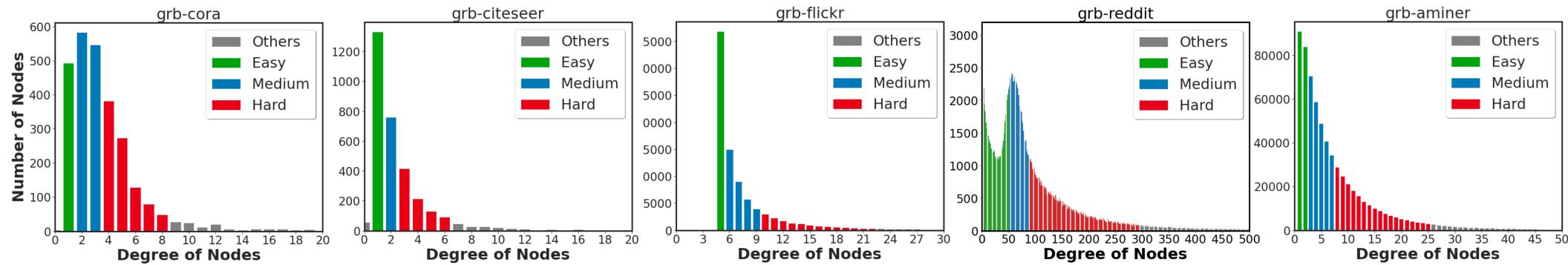
Altogether, GRB serves as a *scalable, unified, modular, and reproducible* benchmark.

# GRB Datasets

- Scalability:** From small to large graphs with hundreds of thousands of nodes (millions of edges).

Dataset	Scale	#Nodes	#Edges	#Feat.	#Classes	Feat. Range	Feat. Range (normalized)
<i>grb-cora</i>	Small	2,680	5,148	302	7	[-2.30, 2.40]	[-0.94, 0.94]
<i>grb-citeseer</i>	Small	3,191	4,172	768	6	[-4.55, 1.67]	[-0.96, 0.89]
<i>grb-flickr</i>	Medium	89,250	449,878	500	7	[-0.90, 269.96]	[-0.47, 1.00]
<i>grb-reddit</i>	Large	232,965	11,606,919	602	41	[-28.19, 120.96]	[-0.98, 0.99]
<i>grb-aminer</i>	Large	659,574	2,878,577	100	18	[-1.74, 1.62]	[-0.93, 0.93]

- Robust-specific splitting:** Various difficulties according to the *average degree* of target nodes.



- Feature normalization:** Unify values to the same scale,  $\mathcal{F} = \frac{2}{\pi} \arctan\left(\frac{\mathcal{F}-\text{mean}(\mathcal{F})}{\text{std}(\mathcal{F})}\right)$

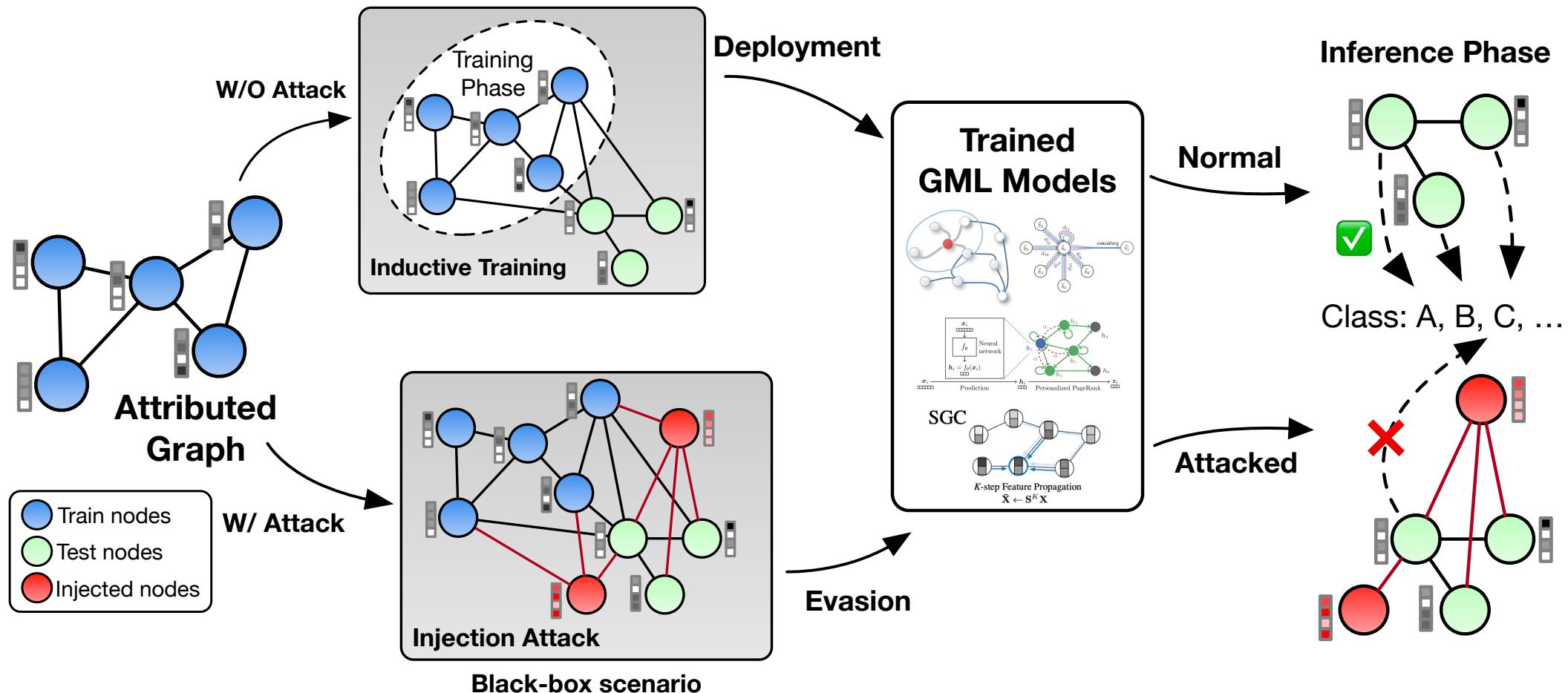
1. X. Zou, Q. Zheng, Y. Dong, X. Guan, E. Kharlamov, J. Lu, and J. Tang. TDGIA: Effective Injection Attacks on Graph Neural Networks. **KDD'21**.

2. Zheng et al., Graph Robustness Benchmark: Benchmarking the Adversarial Robustness of Graph Machine Learning, **NeurIPS'21**.

Code & data for Grand: <https://cogdl.ai/grb/home> <https://github.com/THUDM/grb>

# GRB Attack vs. Defense Scenario

- Example of *graph injection* scenario (*black-box, inductive, evasion*)



**Constraints:** limited number of injected nodes each with limited edges + constrained feature range.

# GRB Modular Coding Framework

GRB Modular Coding Framework		
<b>Dataset</b>	Dataloader	Dataset
	Preprocessing	Splitting Feature normalization
<b>Model</b>	GML Models	GCN, SAGE, GAT, GIN, APPNP, ...
<b>Attack</b>	Modification	RND, DICE, FLIP, STACK, ...
	Injection	RND, FGSM, SPEIT, TDGIA, ...
<b>Defense</b>	Preprocess-based	GNN-SVD, GNNGuard, ...
	Model-based	RobustGNN, Adversarial training, ...
<b>Evaluator</b>	General Metrics	Avg. 3-Min/Max Acc., Weighed Acc., ...
<b>Others</b>	Pipeline	AutoML Visualizer ...
<b>Backend</b>	Pytorch	CogDL DGL

The screenshot shows a Jupyter Notebook interface with two open files: 'training\_example.ipynb' and 'attack\_example.ipynb'. The notebook has a header bar with tabs, file icons, and Python 3 selected. The main area contains sections and code cells.

## 1. Example of using GRB to conduct attacks

### 1.1. Load Dataset

```
[*]: 1 import torch, grb
2 from grb.dataset import Dataset
3
4 dataset = Dataset(name="grb-cora",
5                    data_dir="../../data/",
6                    mode='full',
7                    feat_norm='arctan')
```

Execution started at 2021-11-01 16:22:45

### 1.2. Graph Injection Attack

#### 1.2.1 Train surrogate model

```
[ ]: 1 from grb.model.torch import GCN
2
3 model_sur = GCN(in_features=dataset.num_features,
4                  out_features=dataset.num_classes,
5                  hidden_features=64,
6                  n_layers=2)
```

Example of using GRB

# Example of GRB Leaderboard

*grb-aminer* leaderboard (Top 5 ATK. vs. Top 10 DEF.) in *graph injection* scenario

Attack	Defenses	1	2	3	4	5	6	7	8	9	10	Avg.	Avg. 3-Max	Weighted	
		GAT <sub>ATK</sub>	R-GCN <sub>ATK</sub>	SGCN <sub>LN</sub>	R-GCN	GCN <sub>LN</sub>	GAT <sub>LN</sub>	GIN <sub>LN</sub>	TAGCN <sub>LN</sub>	TAGCN <sub>ATK</sub>	GAT	Acc.	Acc.	Acc.	
1	<b>TDGIA</b>	F	67.69 $\pm$ 0.03	63.62 $\pm$ 0.32	62.20 $\pm$ 0.15	61.99 $\pm$ 0.22	60.38 $\pm$ 1.46	59.69 $\pm$ 1.57	59.59 $\pm$ 0.42	59.06 $\pm$ 1.75	57.24 $\pm$ 5.04	56.63 $\pm$ 6.75	60.81 $\pm$ 1.71	64.52 $\pm$ 2.32	65.74 $\pm$ 0.21
2	<b>SPEIT</b>	F	68.04 $\pm$ 0.03	64.05 $\pm$ 0.04	64.84 $\pm$ 0.04	64.06 $\pm$ 0.04	65.51 $\pm$ 0.02	64.02 $\pm$ 0.04	63.11 $\pm$ 0.02	62.59 $\pm$ 0.04	63.77 $\pm$ 0.06	63.58 $\pm$ 0.06	64.36 $\pm$ 0.02	66.13 $\pm$ 1.38	66.89 $\pm$ 0.02
3	<b>RND</b>	F	67.72 $\pm$ 0.04	64.98 $\pm$ 0.02	65.31 $\pm$ 0.04	64.45 $\pm$ 0.04	66.17 $\pm$ 0.02	67.54 $\pm$ 0.04	64.36 $\pm$ 0.06	64.33 $\pm$ 0.03	66.42 $\pm$ 0.03	66.23 $\pm$ 0.04	65.75 $\pm$ 0.02	67.23 $\pm$ 0.58	67.34 $\pm$ 0.03
4	<b>PGD</b>	F	68.01 $\pm$ 0.02	65.41 $\pm$ 0.01	65.54 $\pm$ 0.03	65.05 $\pm$ 0.03	66.22 $\pm$ 0.02	66.49 $\pm$ 0.04	64.63 $\pm$ 0.04	64.82 $\pm$ 0.04	66.32 $\pm$ 0.02	66.14 $\pm$ 0.04	65.86 $\pm$ 0.01	66.94 $\pm$ 0.76	67.37 $\pm$ 0.02
5	<b>FGSM</b>	F	68.00 $\pm$ 0.02	65.41 $\pm$ 0.02	65.54 $\pm$ 0.04	65.05 $\pm$ 0.04	66.22 $\pm$ 0.02	66.50 $\pm$ 0.06	64.65 $\pm$ 0.04	64.82 $\pm$ 0.03	66.34 $\pm$ 0.03	66.15 $\pm$ 0.06	65.87 $\pm$ 0.01	66.95 $\pm$ 0.75	67.37 $\pm$ 0.01
<b>W/O Attack</b>		F	67.93 $\pm$ 0.00	65.76 $\pm$ 0.00	66.68 $\pm$ 0.00	65.85 $\pm$ 0.00	66.20 $\pm$ 0.00	68.47 $\pm$ 0.00	65.59 $\pm$ 0.00	64.91 $\pm$ 0.00	67.08 $\pm$ 0.00	68.02 $\pm$ 0.00	66.65 $\pm$ 0.00	68.14 $\pm$ 0.24	68.11 $\pm$ 0.00
<b>Avg. Acc.</b>		F	67.90 $\pm$ 0.01	64.87 $\pm$ 0.05	65.02 $\pm$ 0.03	64.41 $\pm$ 0.04	65.12 $\pm$ 0.25	65.45 $\pm$ 0.26	63.65 $\pm$ 0.07	63.42 $\pm$ 0.29	64.53 $\pm$ 0.84	64.46 $\pm$ 1.13	-	-	-
<b>Avg. 3-Min</b>		F	67.78 $\pm$ 0.02	64.22 $\pm$ 0.11	64.12 $\pm$ 0.06	63.50 $\pm$ 0.08	64.02 $\pm$ 0.49	63.39 $\pm$ 0.53	62.35 $\pm$ 0.14	61.99 $\pm$ 0.58	62.44 $\pm$ 1.69	62.11 $\pm$ 2.26	-	-	-
<b>Weighted Acc.</b>		F	67.73 $\pm$ 0.03	63.96 $\pm$ 0.21	63.19 $\pm$ 0.10	62.80 $\pm$ 0.15	62.18 $\pm$ 0.98	61.58 $\pm$ 1.05	61.00 $\pm$ 0.28	60.54 $\pm$ 1.18	59.82 $\pm$ 3.38	59.37 $\pm$ 4.53	-	-	-

- Metric for Attacks:** Avg. Acc., Avg. 3-Max Acc., Weighted Acc.
- Metric for Defenses:** Avg. Acc., Avg. 3-Min Acc., Weighted Acc.
- Avg. 3-Min/Max Acc.:** average accuracy of three most effective attacks or three most robust models.
- Weighted Acc.:** attach higher weight for more effective attacks or more robust models.
- Ranking:** Ranked by Weighted Acc., red and blue indicate the best results of attacks/defenses respectively.
- More results for other difficulties, datasets and scenarios can be found on our website.

● GRB focuses on *generally more effective* methods, rather than comparing a single pair of attack/defense.

1. X. Zou, Q. Zheng, Y. Dong, X. Guan, E. Kharlamov, J. Lu, and J. Tang. TDGIA: Effective Injection Attacks on Graph Neural Networks. **KDD'21**.

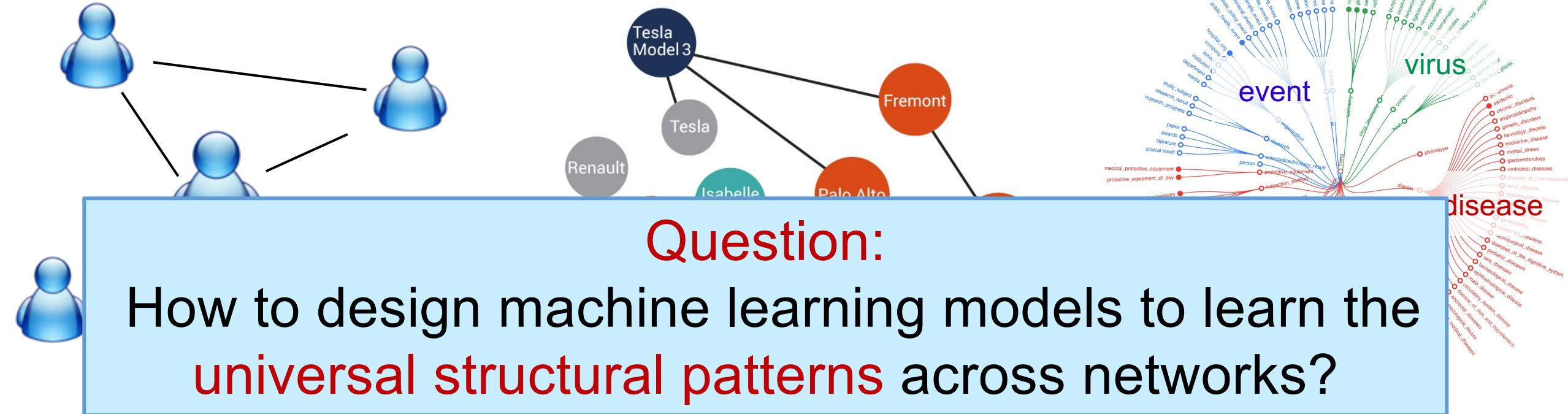
2. Zheng et al., Graph Robustness Benchmark: Benchmarking the Adversarial Robustness of Graph Machine Learning, **NeurIPS'21**.

Code & data for Grand: <https://cogdl.ai/grb/home> <https://github.com/THUDM/grb>



# GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training

# Networked Data



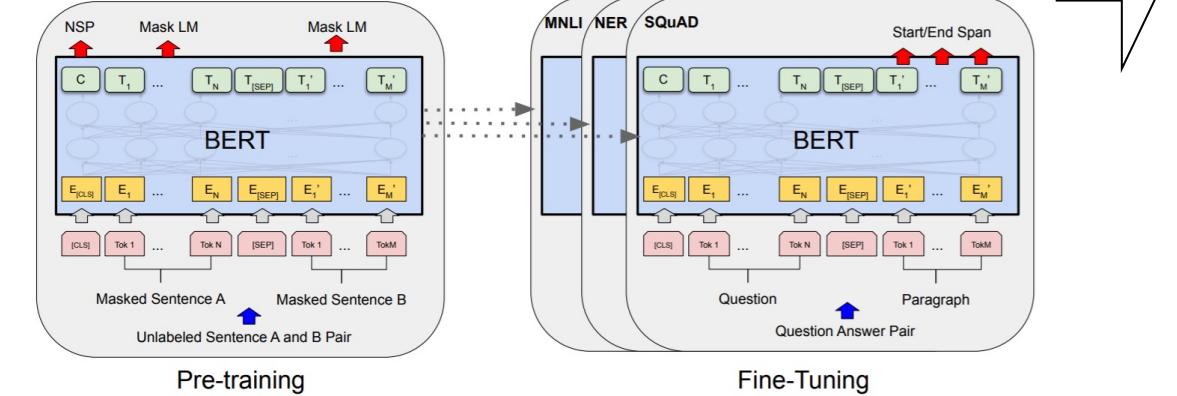
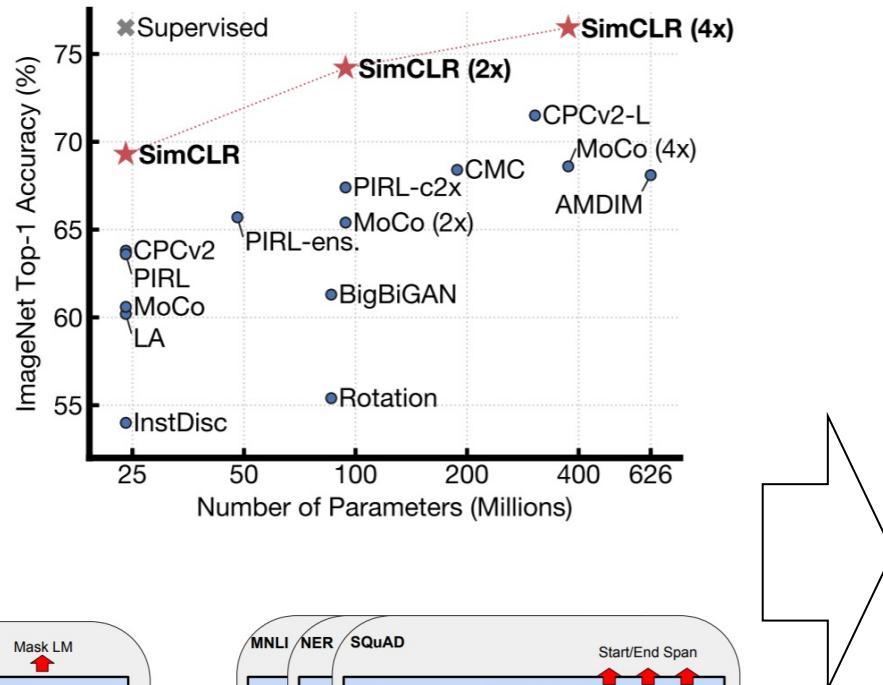
Social Network

Knowledge Graph

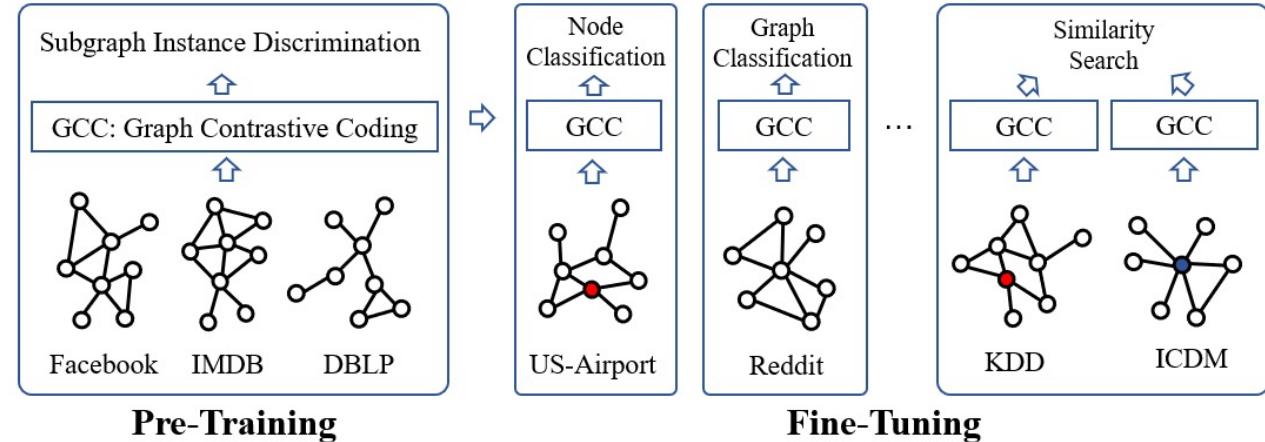
COVID Graph

# Pre-training and Fine-tuning

CV: MoCo,  
SimCLR



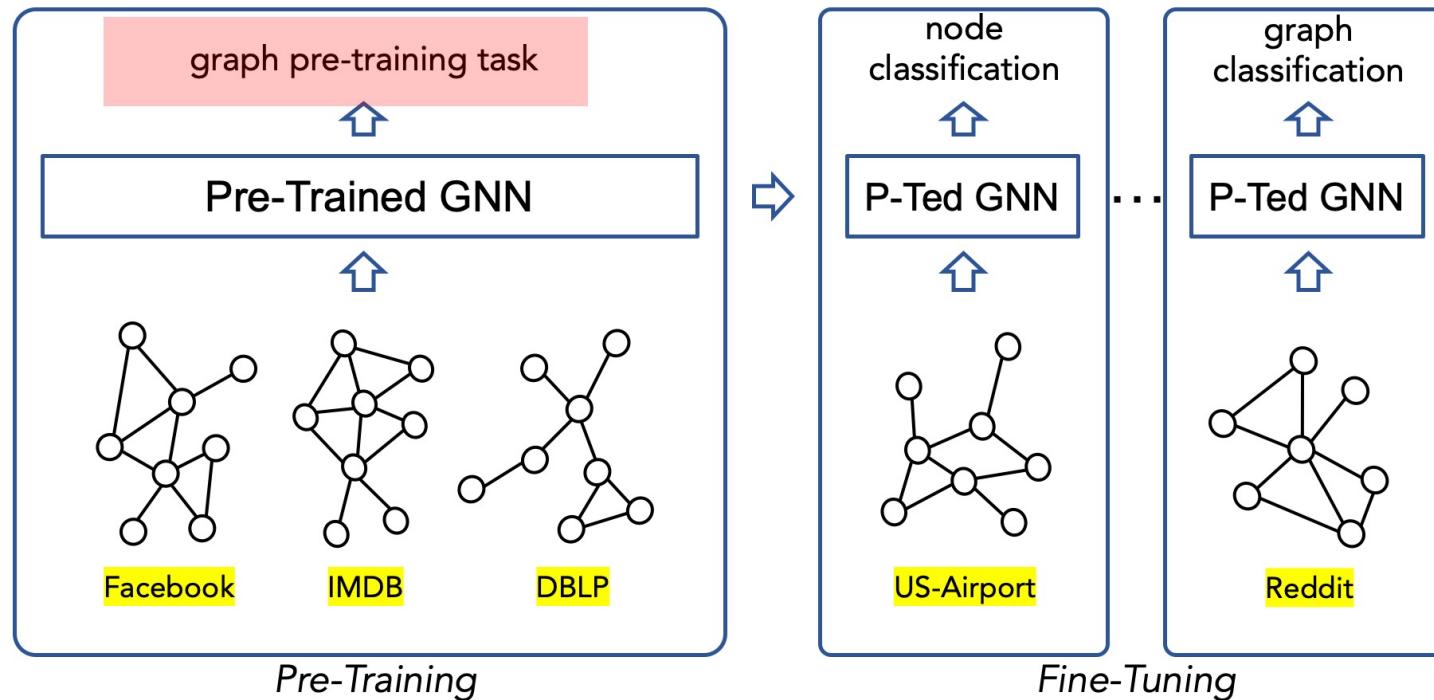
NLP: BERT



Graph Learning  
GCC

# GNN Pre-Training

- Graph pre-training setting:
  - To pre-train from **some graphs**
  - To fine-tune for unseen tasks on **unseen graphs**



- How to do this?
  - Model level: GNNs?
  - Pre-training task: **self-supervised** tasks **across** graphs?

# GNN Pre-Training across Networks

- What are the requirements?
  - **structural similarity**, it maps vertices with similar local network topologies close to each other in the vector space
  - **transferability**, it is compatible with vertices and graphs unseen by the pre-training algorithm

# GNN Pre-Training across Networks

- The Idea: Contrastive learning
  - **pre-training task:** instance discrimination
  - **InfoNCE objective:** output instance representations that are capable of capturing the similarities between instances

$$\mathcal{L} = -\log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{i=0}^K \exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}$$

- query instance  $x^q$
- query  $\mathbf{q}$  (embedding of  $x^q$ ), i.e.,  $\mathbf{q} = f(x^q)$
- dictionary of keys  $\{\mathbf{k}_0, \mathbf{k}_1, \dots, \mathbf{k}_K\}$
- key  $\mathbf{k} = f(x^k)$

- Contrastive learning for graphs?

- Q1: How to define instances in graphs?
- Q2: How to define (dis) similar instance pairs in and across graphs?
- Q3: What are the proper graph encoders?

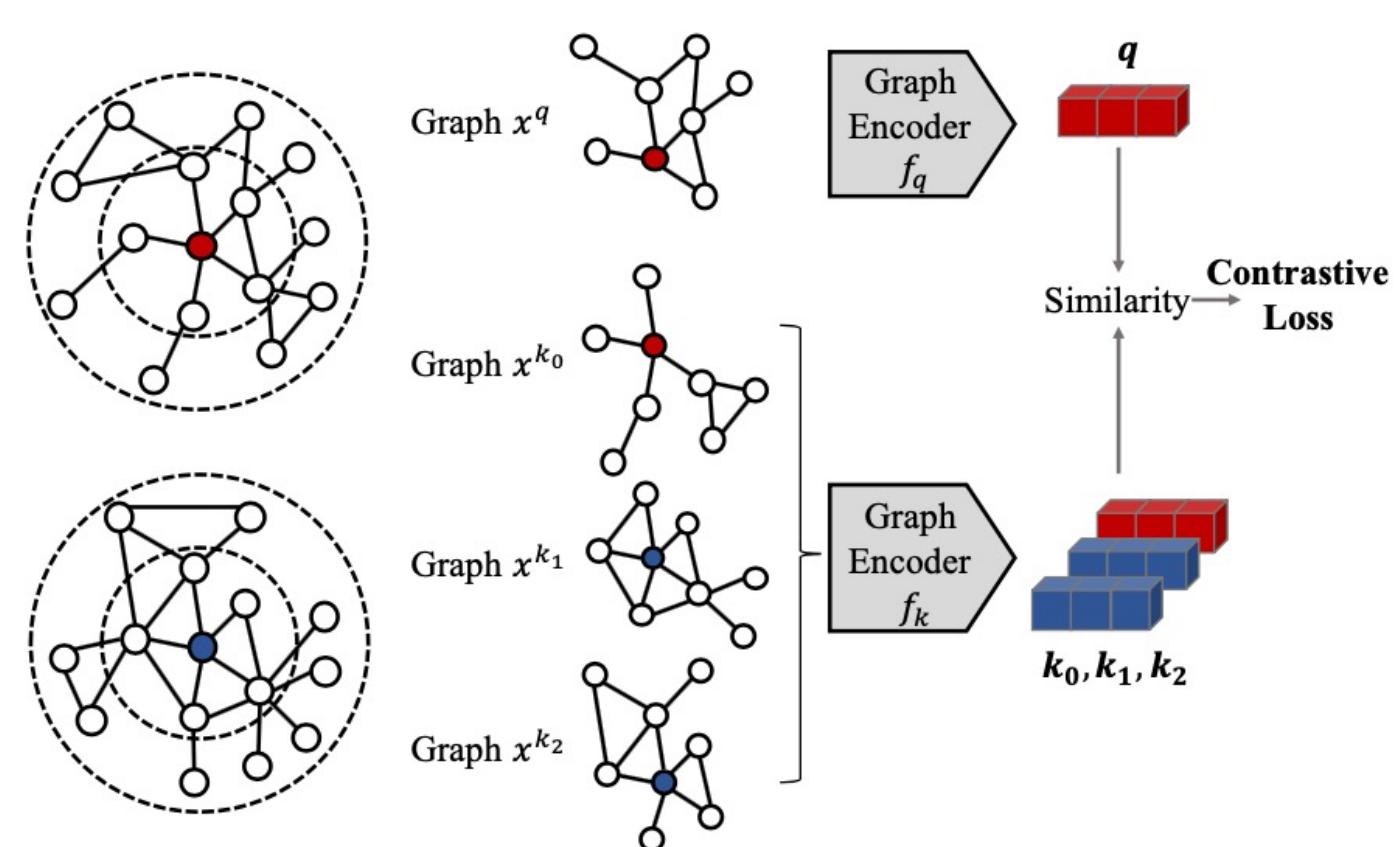
# Graph Contrastive Coding (GCC)

- Contrastive learning for graphs

- Q1: How to define instances in graphs?
- Q2: How to define (dis) similar instance pairs in and across graphs?
- Q3: What are the proper graph representations?

$$\mathcal{L} = -\log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{i=0}^K \exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}$$

Subgraph instance discrimination



# GCC Pre-Training / Fine-Tuning

- pre-train on six graphs

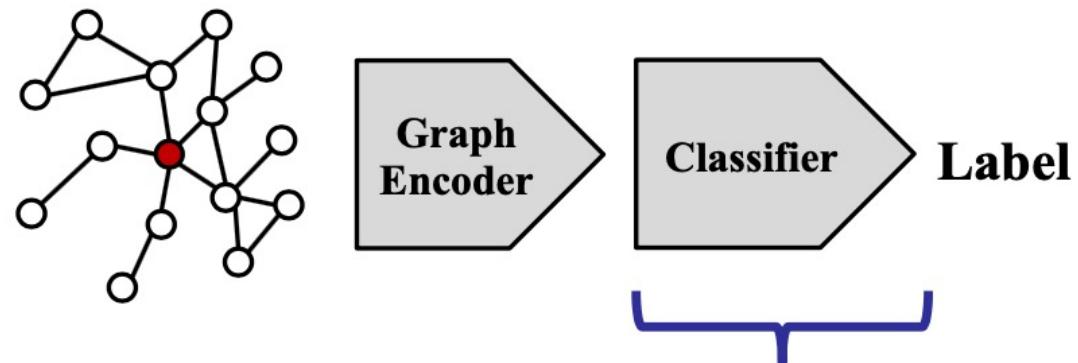
Code & Data for GCC:  
<https://github.com/THUDM/GCC>

Dataset	Academia	DBLP (SNAP)	DBLP (NetRep)	IMDB	Facebook	LiveJournal
$ V $	137,969	317,080	540,486	896,305	3,097,165	4,843,953
$ E $	739,384	2,099,732	30,491,458	7,564,894	47,334,788	85,691,368

- fine-tune on **different** graphs

- US-Airport & AMiner academic graph
  - Node classification
- COLLAB, RDT-B, RDT-M, & IMDB-B/M
  - Graph classification
- AMiner academic graph
  - Similarity search

**Full Fine-tuning**

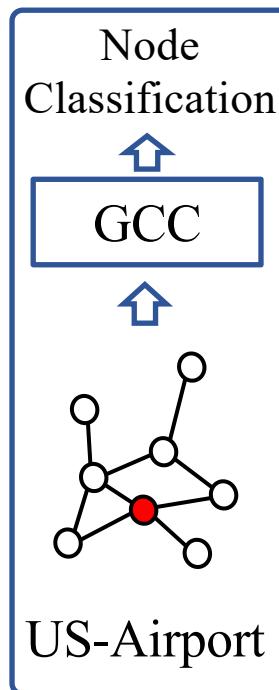


- The base GNN
  - Graph Isomorphism Network (GIN)

**Freezing Fine-tuning**

# Result 1: Node Classification

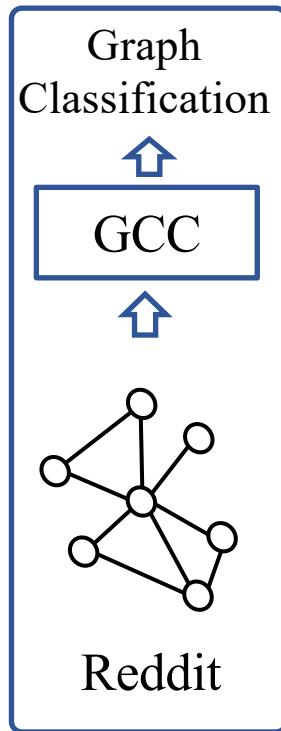
- Setup
  - US-Airport
  - AMiner academic graph



Datasets	US-Airport	H-index
$ V $	1,190	5,000
$ E $	13,599	44,020
ProNE	62.3	69.1
GraphWave	60.2	70.3
Struc2vec	<b>66.2</b>	> 1 Day
GCC (E2E, freeze)	64.8	<b>78.3</b>
GCC (MoCo, freeze)	65.6	75.2
GCC (rand, full)	64.2	76.9
GCC (E2E, full)	<b>68.3</b>	80.5
GCC (MoCo, full)	67.2	<b>80.6</b>

# Result 2: Graph Classification

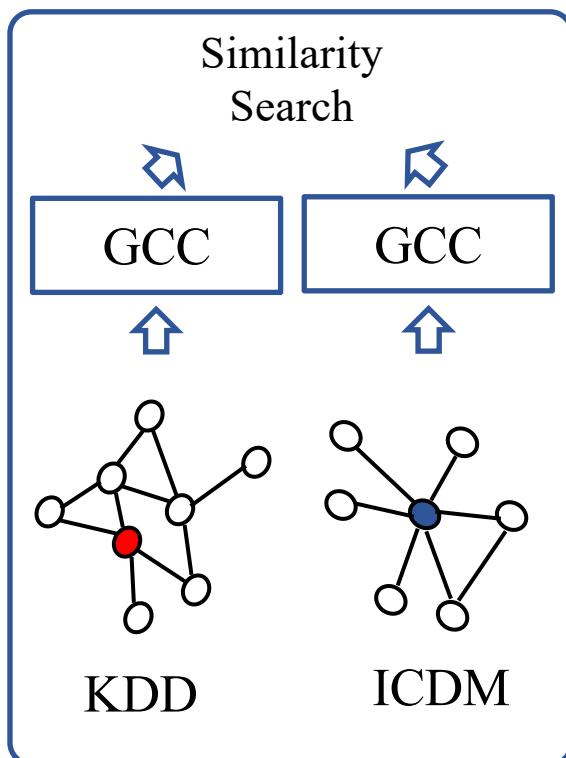
- Setup
  - COLLAB, RDT-B, RDT-M, & IMDB-B, IMDB-M



Datasets	IMDB-B	IMDB-M	COLLAB	RDT-B	RDT-M
# graphs	1,000	1,500	5,000	2,000	5,000
# classes	2	3	3	2	5
Avg. # nodes	19.8	13.0	74.5	429.6	508.5
DGK	67.0	44.6	73.1	78.0	41.3
graph2vec	71.1	50.4	–	75.8	47.9
InfoGraph	<b>73.0</b>	<b>49.7</b>	–	82.5	53.5
GCC (E2E, freeze)	71.7	49.3	74.7	87.5	52.6
GCC (MoCo, freeze)	72.0	49.4	<b>78.9</b>	<b>89.8</b>	<b>53.7</b>
DGCNN	70.0	47.8	73.7	–	–
GIN	<b>75.6</b>	<b>51.5</b>	80.2	<b>89.4</b>	<b>54.5</b>
GCC (rand, full)	<b>75.6</b>	50.9	79.4	87.8	52.1
GCC (E2E, full)	70.8	48.5	79.0	86.4	47.4
GCC (MoCo, full)	73.8	50.3	<b>81.1</b>	87.6	53.0

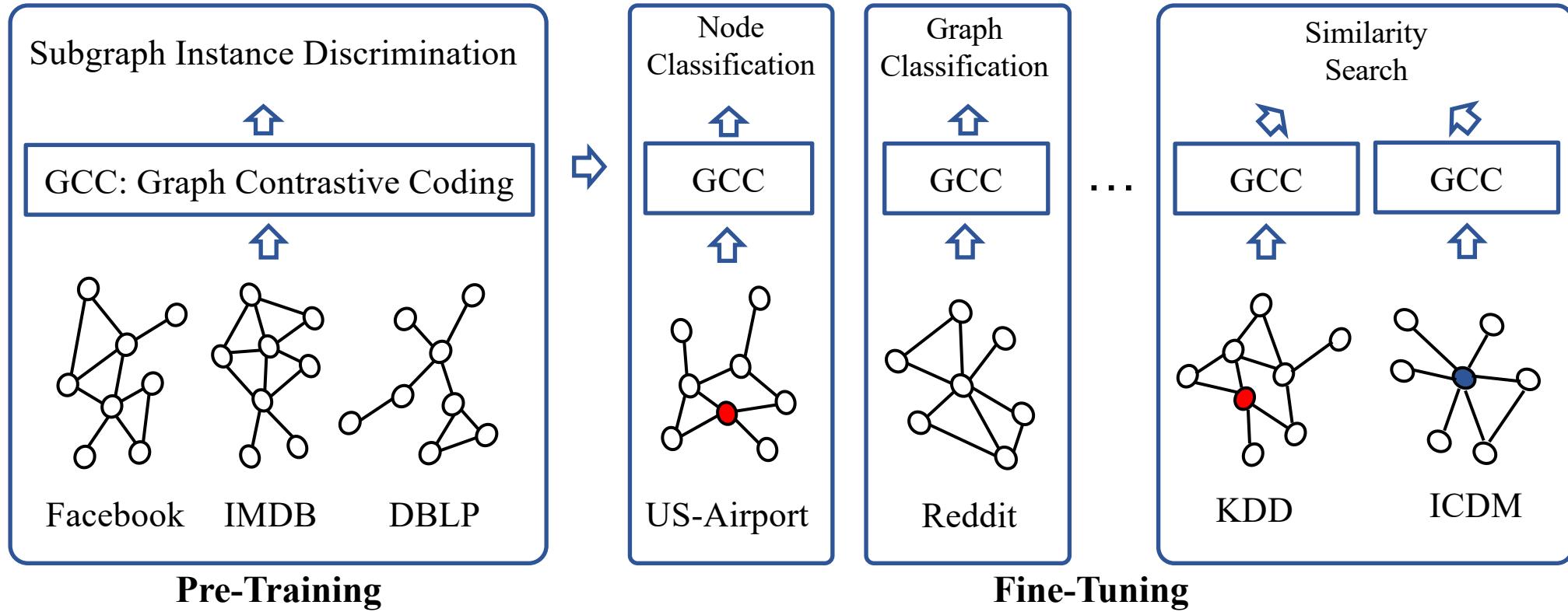
# Result 3: Top-k Similarity Search

- Setup
  - AMiner academic graph



	KDD-ICDM		SIGIR-CIKM		SIGMOD-ICDE	
$ V $	2,867	2,607	2,851	3,548	2,616	2,559
$ E $	7,637	4,774	6,354	7,076	8,304	6,668
# group truth		697		874		898
$k$	20	40	20	40	20	40
Random	0.0198	0.0566	0.0223	0.0447	0.0221	0.0521
RoLX	0.0779	0.1288	0.0548	0.0984	0.0776	0.1309
Panther++	0.0892	0.1558	<b>0.0782</b>	0.1185	0.0921	0.1320
GraphWave	0.0846	<b>0.1693</b>	0.0549	0.0995	<b>0.0947</b>	<b>0.1470</b>
GCC (E2E)	<b>0.1047</b>	0.1564	0.0549	<b>0.1247</b>	0.0835	0.1336
GCC (MoCo)	0.0904	0.1521	0.0652	0.1178	0.0846	0.1425

# GNN Pre-Training



1. Jiezhong Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. **KDD** 2020.

2. Code & Data for GCC: <https://github.com/THUDM/GCC>

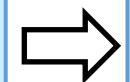
# GRL: NE&GNN

Network  
Embedding

Matrix  
Factorization

Graph Neural  
Networks

GNN  
Pre-Training



Learning GNNs with **CogDL**



<https://alchemy.tencent.com/>



**CogDL**

<https://github.com/thudm/cogdl>

It is **hard** to apply GNNs to real-world applications...

- **Billion-scale** real-world graphs!
  - How to store large-scale graphs, GPU memory bounded
  - Large cost of model training
- Potential **issues** in training GNNs
  - Over-fitting, Over-smoothing
- Methodologies
  - “Millions of” papers...
  - Which one really works? Reproducibility...

# CogDL

**CogDL: An Extensive Research Toolkit for**

$$h_v^k \leftarrow \sigma \left( W \cdot \text{MEAN} \left( \{h_v^{k-1}\} \cup \{h_u^{k-1}, \forall u \in \mathcal{N}(v)\} \right) \right) \quad \log \left( \frac{\text{Vol}(G)}{b} D^{-1} A D^{-1} \right)$$

**Deep Learning on Graphs**

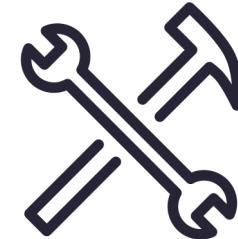
KEG, Tsinghua University

# CogDL Introduction

Vision

CogDL aims at providing researchers and developers with easy-to-use APIs, reproducible results, and high efficiency for most graph tasks and applications.

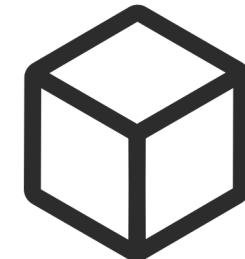
Philosophy



Easy-to-use

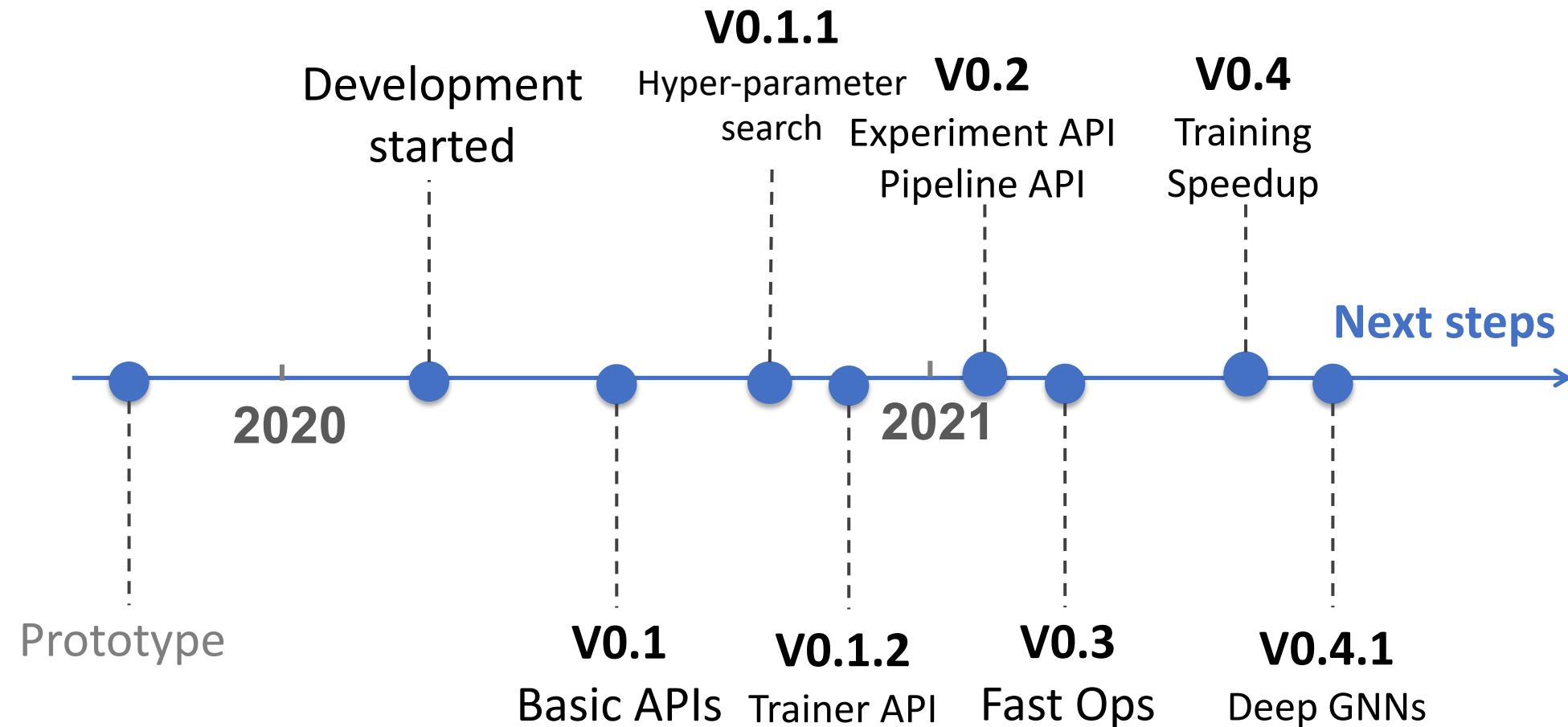


Reproducibility



Efficiency

# CogDL Development

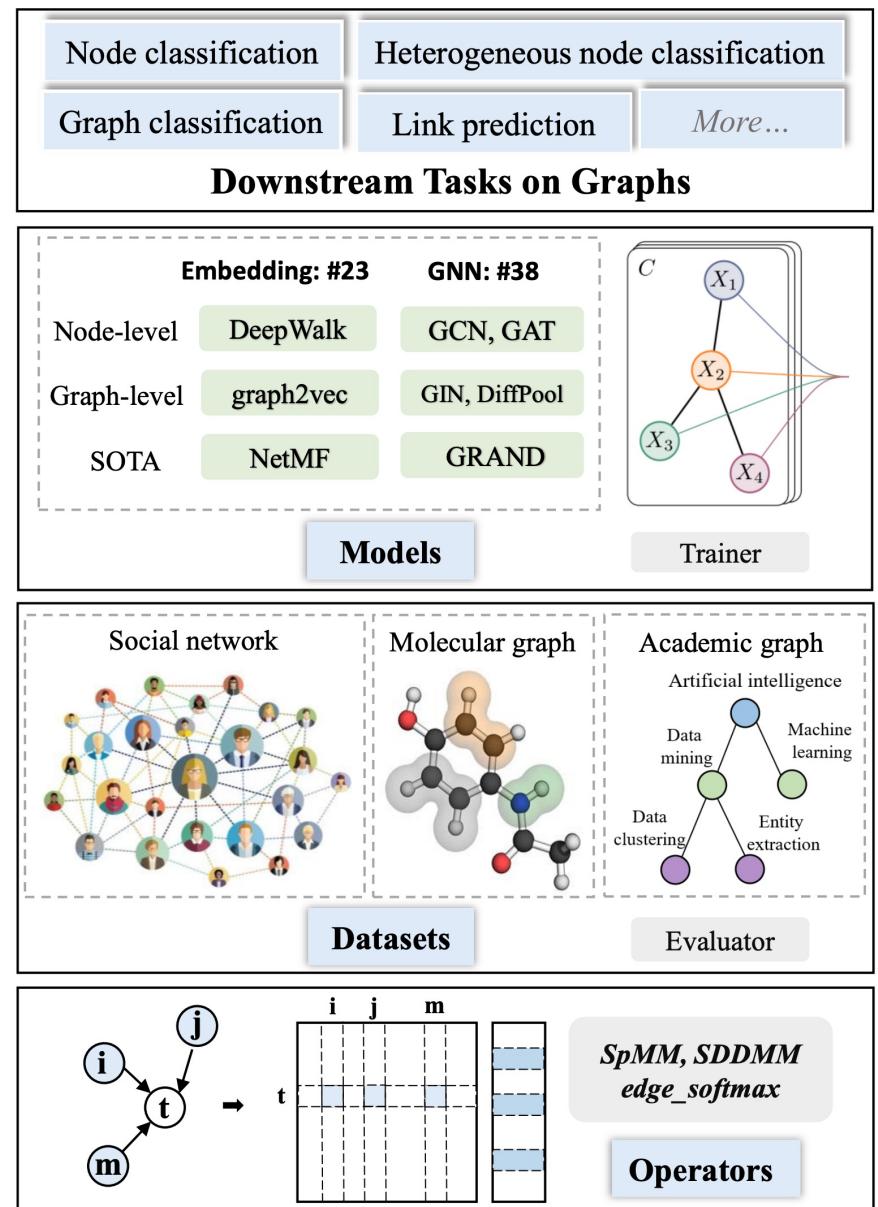


Prerequisite: PyTorch environment  
CogDL installation: **pip install cogdl**  
or git clone <https://github.com/THUDM/cogdl>

# Tasks, Datasets, Models in CogDL

- >10 Tasks
  - Node/graph classification, etc.
- >60 Datasets
  - Social networks, academics, molecular
- >70 models

Task	Characteristics	Models
Node Classification (Unsupervised)	<i>MF based</i>	SpectralClustering [7]
	<i>With high-order neighborhood</i>	NetMF [1], ProNE [2], NetSMF [3], HOPE [8], GraRep [9]
	<i>Skip-gram</i>	LINE [5]
Node Classification (with GNN)	<i>Semi-supervised</i>	GCN [18], GAT [16], JK-Net [20], ChebyNet [21], GCNII [11]
	<i>With random propagation</i>	DropEdge [17], Graph-Unet [14], GraphSAINT [39], GRAND [10]
	<i>With diffusion</i>	GDC [15], APPNP [13], GRAND [10], PPRGo [40]
Graph Classification	<i>Self-supervised</i>	Contrastive methods
		MVGRL [12], DGI [19]
	<i>Supervised</i>	CNN method
		PATCHY_SAN [41]
		Global pooling
		GIN [42], SortPool [43], DGCNN [44]
		Hierarchical pooling
	<i>Unsupervised</i>	Kernel methods
		DGK [47], graph2vec [48]
		GNN method
		Infograph [49]



# Experiment API

- Feed task, dataset, model, (hyper-parameters), (search space)

```
from cogdl import experiment

# basic usage
experiment(task="node_classification", dataset="cora", model="gcn")

# set other hyper-parameters
experiment(task="node_classification", dataset="cora", model="gcn", hidden_size=32, max_epoch=200)

# run over multiple models on different seeds
experiment(task="node_classification", dataset="cora", model=["gcn", "gat"], seed=[1, 2])

# automl usage
def func_search(trial):
    return {
        "lr": trial.suggest_categorical("lr", [1e-3, 5e-3, 1e-2]),
        "hidden_size": trial.suggest_categorical("hidden_size", [32, 64, 128]),
        "dropout": trial.suggest_uniform("dropout", 0.5, 0.8),
    }

experiment(task="node_classification", dataset="cora", model="gcn", seed=[1, 2], func_search=func_search)
```

# Results of Experiment API

```
In [3]: from cogdl import experiment

# basic usage
experiment(task="node_classification", dataset="cora", model="gcn")

Failed to load C version of sampling, use python version instead.
Namespace(activation='relu', checkpoint=None, cpu=False, dataset='cora', device_id=[0], dropout=0.5, fast_spmm=False, hidden_size=64, inference=False, lr=0.01, max_epoch=500, missing_rate=0, model='gcn', norm=None, num_classes=None, num_features=None, num_layers=2, patience=100, residual=False, save_dir='.', save_model=None, seed=1, task='node_classification', trainer=None, use_best_config=False, weight_decay=0.0005)
Downloading https://cloud.tsinghua.edu.cn/d/6808093f7f8042bfaf0/files/?p=%2Fcora.zip&dl=1
unpacking cora.zip
Processing...
Done!

Epoch: 455, Train: 1.0000, Val: 0.7920, ValLoss: 0.7300: 89%|██████████| 446/500 [00:03<00:00, 120.41it/s]
Valid accuracy = 0.7940
Test accuracy = 0.8100
| Variant      | Acc          | ValAcc        |
|---|---|---|
| ('cora', 'gcn') | 0.8100±0.0000 | 0.7940±0.0000 |
```

```
Out[3]: defaultdict(list, {('cora', 'gcn'): [{'Acc': 0.81, 'ValAcc': 0.794}]}))
```

# Results of Node Classification

- Two kinds of models:
  - Semi-supervised: GCN, GAT, GRAND, ...
  - Self-supervised: MVGRL, DGI
- Citation networks: Cora, Citeseer, Pubmed

Rank	Method	Cora	Citeseer	Pubmed	Reproducible
1	GRAND [12]	84.8	75.1	82.4	Yes
2	GCNII [7]	85.1	71.3	80.2	Yes
3	MVGRL [20]	83.6 ↓	73.0	80.1	Partial
4	APPNP [26]	<b>84.3 ↑</b>	72.0	80.0	Yes
5	Graph-Unet [15]	83.3 ↓	71.2 ↓	79.0	Partial
6	GDC [27]	82.5	72.1	79.8	Yes
7	GAT [53]	82.9	71.0	78.9	Yes
8	DropEdge [38]	82.1	72.1	79.7	Yes
9	GCN [25]	<b>82.3 ↑</b>	<b>71.4 ↑</b>	79.5	Yes
10	DGI [52]	82.0	71.2	76.5	Yes
11	JK-net [58]	81.8	69.5	77.7	Yes
12	Chebyshev [8]	79.0	69.8	68.6	Yes

# Results of Graph Classification

- Two kinds of models
  - Self-supervised: InfoGraph, graph2vec, DGK
  - Supervised: GIN, DiffPool, SortPool, ...
- Two types of graphs
  - Bioinformatics: MUTAG, PTC, NCI1, PROTEINS
  - Social networks: IMDB-B/M, COLLAB, REDDIT-B

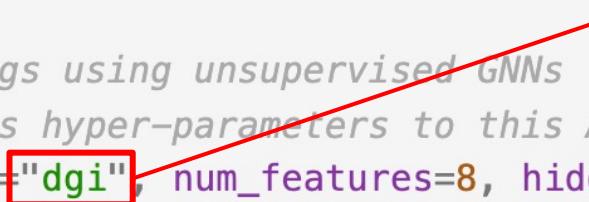
Algorithm	MUTAG	PTC	NCI1	PROTEINS	IMDB-B	IMDB-M	COLLAB	REDDIT-B	<i>Reproducible</i>
GIN [57]	92.06	67.82	81.66	75.19	76.10	51.80	79.52	83.10 ↓	Yes
InfoGraph [42]	88.95	60.74	76.64	73.93	74.50	51.33	79.40	76.55	Yes
DiffPool [62]	85.18	58.00	69.09	75.30	72.50	50.50	79.27	81.20	Yes
SortPool [67]	87.25	62.04	<b>73.99 ↑</b>	74.48	75.40	50.47	<b>80.07 ↑</b>	78.15	Yes
graph2vec [31]	83.68	54.76 ↓	71.85	73.30	73.90	52.27	<b>85.58 ↑</b>	91.77	Yes
PATCHY_SAN [32]	86.12	61.60	69.82	75.38	<b>76.00 ↑</b>	46.40	74.34	60.61	Yes
DGCNN [56]	83.33	56.72	65.96	66.75	71.60	49.20	77.45	86.20	Yes
SAGPool [28]	71.73 ↓	59.92	72.87	74.03	74.80	51.33	/	89.21	Yes
DGK [59]	85.58	57.28	/	72.59	55.00 ↓	40.40 ↓	/	/	<i>Partial</i>

# Pipeline API

- Feed application, model, (hyper-parameters)

```
import numpy as np
from cogdl import pipeline

# build a pipeline for generating embeddings
# pass model name with its hyper-parameters to this API
generator = pipeline("generate-emb", model="prone")  
  
# generate embedding by an unweighted graph
edge_index = np.array([[0, 1], [0, 2], [0, 3], [1, 2], [2, 3]])
outputs = generator(edge_index)
print(outputs)

# build a pipeline for generating embeddings using unsupervised GNNs
# pass model name and num_features with its hyper-parameters to this API
generator = pipeline("generate-emb", model="dgi", num_features=8, hidden_size=4)  
  
outputs = generator(edge_index, x=np.random.randn(4, 8))
print(outputs)
```

- prone
- netmf
- netsmf
- deepwalk
- line
- node2vec
- hope
- sdne
- grarep
- dngr
- spectral

- unsup\_graphsage
- dgi
- mvgrl
- grace

# All with CogDL

- Efficiency
  - Graph storage in CogDL
  - Sparse operators in CogDL
  - Training on large-scale graphs
  - Training very deep GNNs
- Customization
  - Customized usage in CogDL
- Benchmarks:
  - Self-supervised learning
  - Heterogeneous Graph Benchmark (HGB)
  - Graph Robustness Benchmark (GRB)
- Applications:
  - Recommendation

# Sparse Storage of Adjacency Matrix

- COO format:
  - (row, col) or (row, col, value), size:  $|E|^2/3$
  - $[[0,0,1], [0,2,2], [1,2,3], [2,0,4], [2,1,5], [2,2,6]]$
- CSR format:
  - row\_ptr: size  $|V|+1$
  - col\_indices: size  $|E|$
  - value: size  $|E|$
  - $[0, 2, 3, 6], [0, 2, 2, 0, 1, 2], [1, 2, 3, 4, 5, 6]$

$$\begin{bmatrix} 1 & 0 & 2 \\ 0 & 0 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

# Graph Storage in CogDL

`class Graph: (defined in cogdl.data)`

- `x`: node feature matrix
- `y`: node labels
- `edge_index`: COO format matrix
- `edge_weight`: edge weight (if exists)
- `edge_attr`: edge attributes (if exists)
- `row_ptr`: row index pointer for CSR matrix
- `col_indices`: column indices for CSR matrix

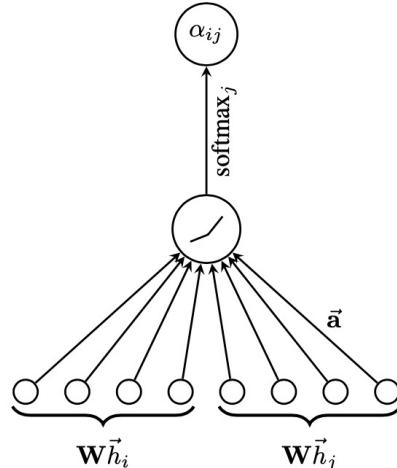
# Usage of CogDL's Graph

- Graph Initialization
  - `g = Graph(edge_index=edge_index)`
  - `g.edge_weight = torch.rand(n)`
- Commonly used operators:
  - `add_self_loops()`
  - `sym_norm()`
  - `degrees()`
  - `subgraph()`
  - ...

# Recall Sparse Operators in GNNs

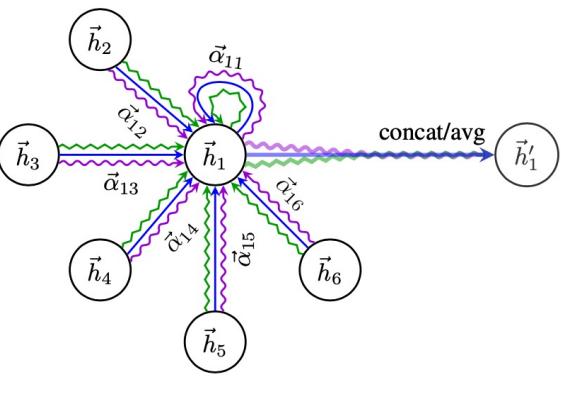
- GCN (Sparse Matrix-Matrix Multiplication, SpMM)

$$\mathbf{H}^{(i+1)} = \mathbf{A}\mathbf{H}^{(i)}\mathbf{W}$$



- GAT (Edge-wise-softmax)

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})}$$



- GAT (Multi-Head SpMM)

$$\mathbf{h}_i = \text{CONCAT} \left( \sigma \left( \sum_{j \in N_i} \alpha_{ij}^k \mathbf{W}^k \mathbf{h}_j \right) \right)$$

# GCN/GAT Layer in CogDL

```
# GCN Layer  
# graph: cogdl.data.Graph  
# x: node features  
# weight: parameters  
  
h = torch.mm(x, weight)  
h = spmm(graph, h)  
out = torch.relu(h)
```

$$H^{(i+1)} = AH^{(i)}W^{(i)}$$

```
# GAT Layer  
# graph: cogdl.data.Graph  
# h: node features  
# h_score: importance score of edge  
  
edge_attention = mul_edge_softmax(graph, edge_score)  
h = mh_spmm(graph, edge_attention, h)  
out = torch.cat(h, dim=1)
```

$$\alpha_{ij} = softmax(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})}$$

$$h_i = CONCAT \left( \sigma \left( \sum_{j \in N_i} \alpha_{ij}^k W^k h_j \right) \right)$$

# Implementation of GCN/GAT Layer

```
class GCNLayer(nn.Module):
    """
    Simple GCN layer, similar to https://arxiv.org/abs/1706.10219
    """

    def __init__(self, in_features, out_features):
        super(GCNLayer, self).__init__()

        self.reset_parameters()

        self.forward = self._forward

    def _forward(self, graph, x):
        support = torch.mm(x, self.weight)
        out = spmm(graph, support)
```

$$H^{(i+1)} = AH^{(i)}W^{(i)}$$

```
class GATLayer(nn.Module):
    """
    Sparse version GAT layer, similar to https://arxiv.org/abs/1710.10903
    """

    def __init__(self, in_features, out_features, nhead=1, alpha=0.2, dropout=0.5):
        super(GATLayer, self).__init__()

        self.reset_parameters()

        self.forward = self._forward

    def _reset_parameters(self):
        glorot(self.weight)
        glorot(self.attention)

    def _forward(self, graph, x):
        h = torch.matmul(x, self.weight)
        edge_score = self.compute_edge_score(graph, x, h)
        # edge_attention: E * H
        edge_attention = mul_edge_softmax(graph, edge_score)
        edge_attention = self.dropout(edge_attention)

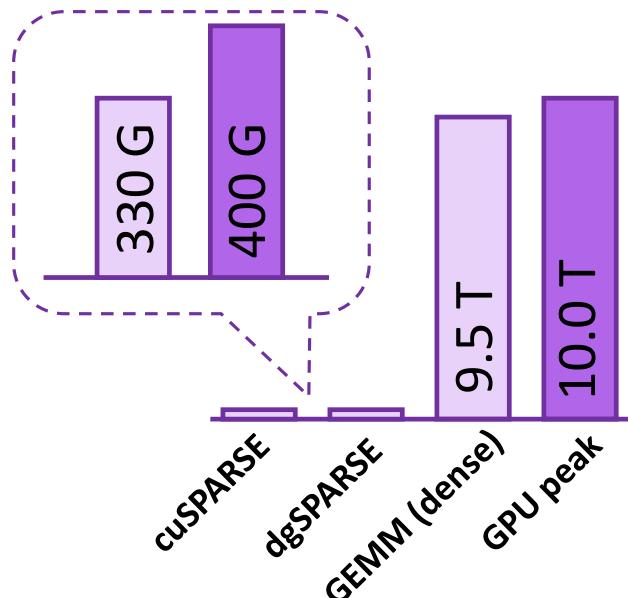
        out = mh_spmm(graph, edge_attention, h)
```

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})}$$

$$h_i = \text{CONCAT} \left( \sigma \left( \sum_{j \in N_i} \alpha_{ij}^k W^k h_j \right) \right)$$

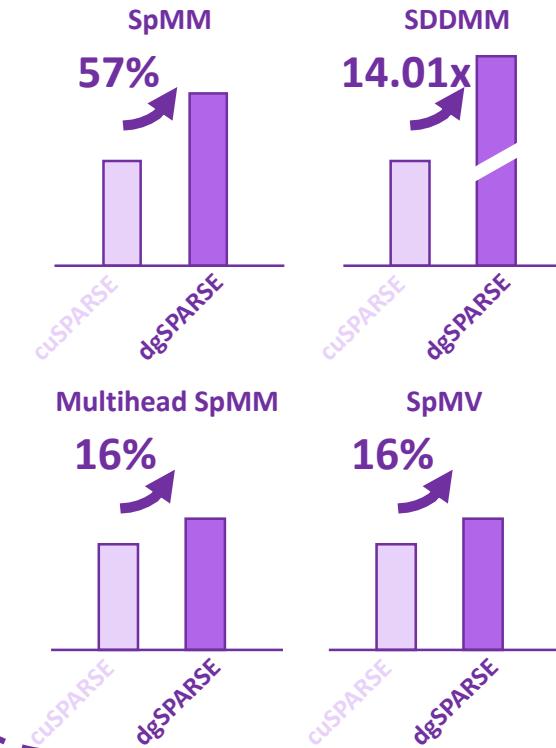
# Efficient Sparse Kernels

Lack efficient sparse kernels

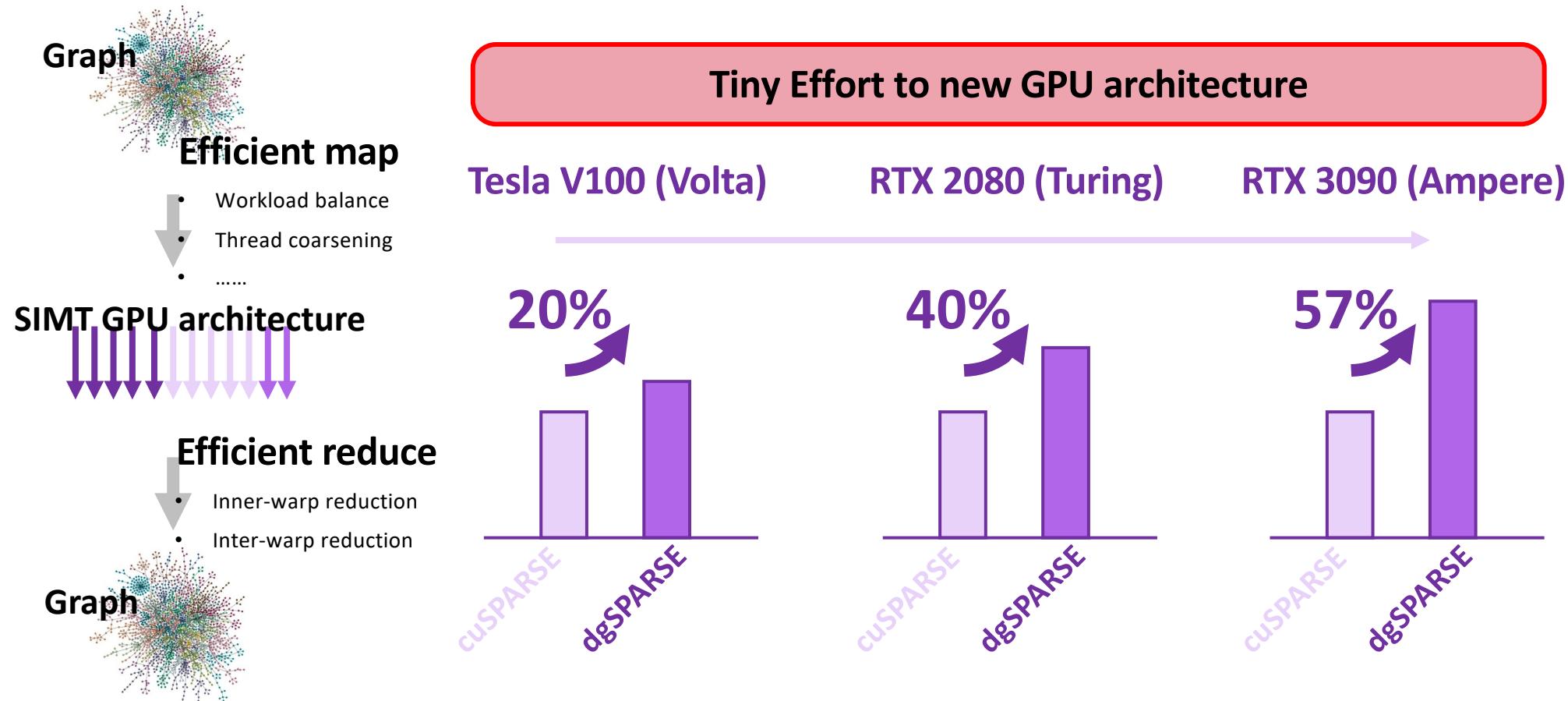


Big gap between FLOPS of sparse kernels  
and hardware peak performance.

dgSPARSE, Deep Graph SPARSE  
Efficient Implementation on GPUs



# Deep Graph Sparse (dgSPARSE) Library



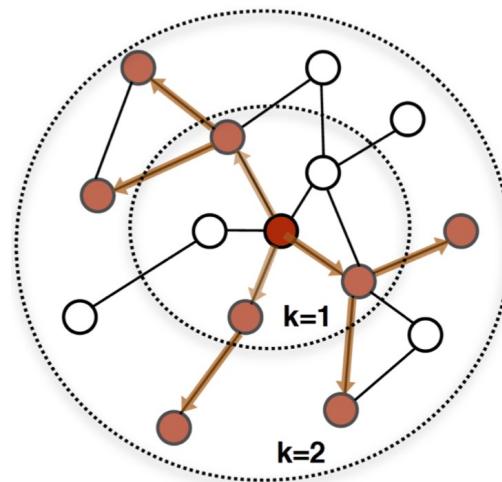
# Performance of GCN/GAT model

- Setting: 2-layer GCN/GAT, hidden size=128
- Supported by dgSPAESE

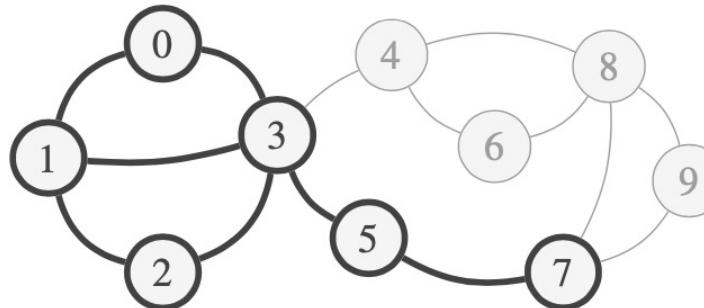
Model	GPU	Dataset	Tranining per epoch (s)				Inference per epoch (s)			
			torch	PyG	DGL	CogDL	torch	PyG	DGL	CogDL
GCN	2080Ti (11G)	Flickr	0.025	0.0084	0.012	0.085	0.012	0.0034	0.007	0.0035
		Reddit	0.445	0.122	0.102	0.081	0.218	0.045	0.049	0.039
		Yelp	0.412	0.151	0.151	0.110	0.191	0.053	0.063	0.040
	3090 (24G)	Flickr	0.017	0.006	0.008	0.007	0.008	0.002	0.004	0.002
		Reddit	0.263	0.062	0.060	0.050	0.127	0.022	0.0314	0.022
		Yelp	0.230	0.081	0.080	0.062	0.106	0.029	0.036	0.023
GAT	2080Ti (11G)	PubMed	0.017	0.016	0.011	0.012	0.004	0.006	0.004	0.003
		Flickr	0.082	0.090	0.047	0.056	0.023	0.030	0.019	0.014
		*Reddit	—‡	—‡	0.406	0.537	—‡	—‡	0.163	0.086
	3090 (24G)	PubMed	0.043	0.011	0.011	0.016	0.004	0.004	0.003	0.002
		Flickr	0.097	0.059	0.033	0.044	0.021	0.024	0.013	0.009
		Reddit	0.671	—‡	0.301	0.373	0.112	—‡	0.113	0.088
		Yelp	0.614	—‡	0.294	0.404	0.118	—‡	0.105	0.113

# Training on Large-scale Graphs

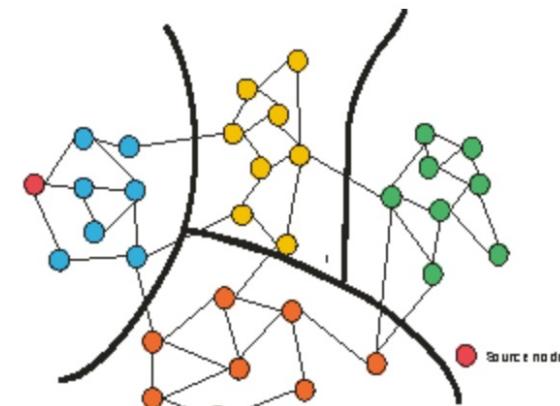
- Billion-scale social networks and recommender systems
- Main challenge: GPU memory bounded!
- Training GNNs via mini-batch sampling



Neighbor Sampling  
(NeurIPS '17)



GraphSAINT  
(ICLR '20)

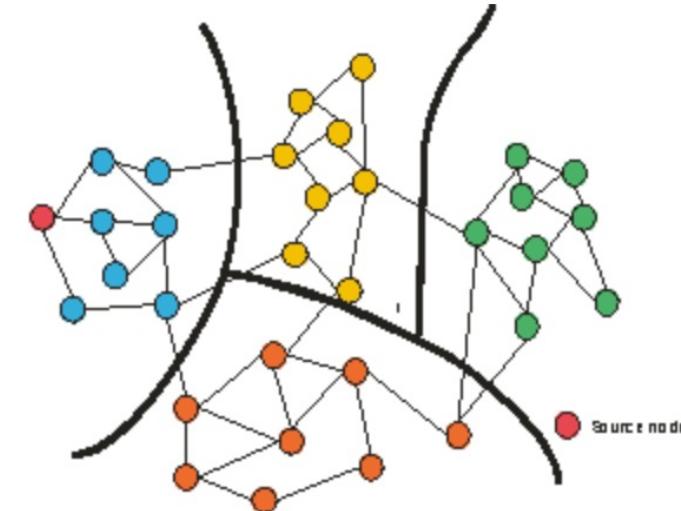


ClusterGCN  
(KDD '19)

1. Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. NeurIPS '17.
2. Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-GCN: An efficient algorithm for training deep and large graph convolutional networks. KDD '19.
3. Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graphsaint: Graph sampling based inductive learning method. In ICLR '20.

# ClusterGCN

- Graph partition (METIS)
- Train GNNs via mini-batch
- Memory:  $O(NFL) \rightarrow O(bFL)$



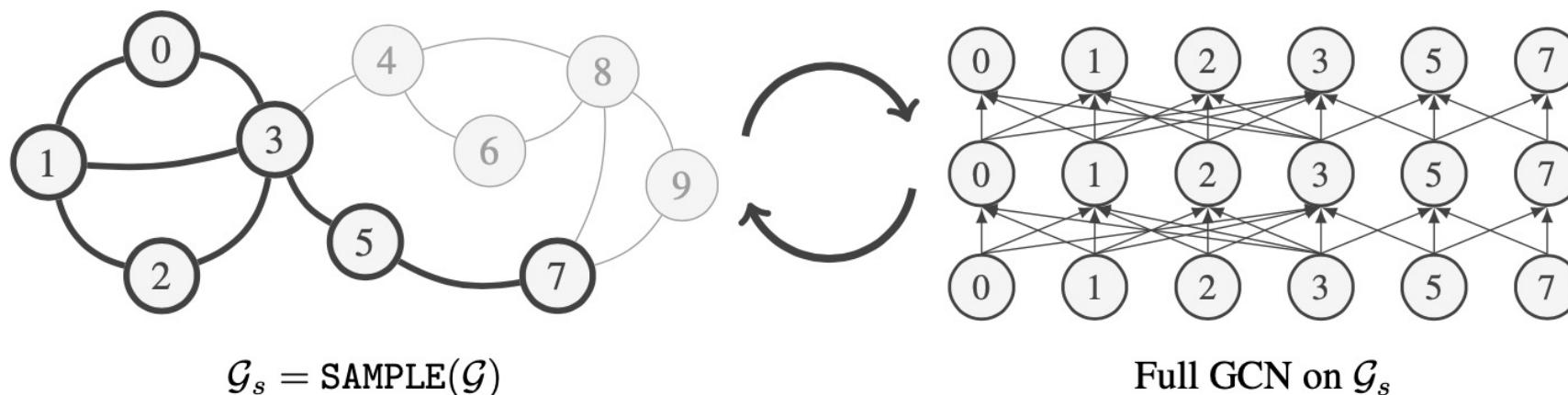
	GCN [9]	Vanilla SGD	GraphSAGE [5]	FastGCN [1]	VR-GCN [2]	Cluster-GCN
Time complexity	$O(L\ A\ _0F + LNF^2)$	$O(d^LNF^2)$	$O(r^LNF^2)$	$O(rLNF^2)$	$O(L\ A\ _0F + LNF^2 + r^LNF^2)$	$O(L\ A\ _0F + LNF^2)$
Memory complexity	$O(LNF + LF^2)$	$O(bd^L F + LF^2)$	$O(br^L F + LF^2)$	$O(brLF + LF^2)$	$O(LNF + LF^2)$	$O(bLF + LF^2)$

	Time		Memory		Test F1 score	
	VRGCN	Cluster-GCN	VRGCN	Cluster-GCN	VRGCN	Cluster-GCN
Amazon2M (2-layer)	337s	1223s	7476 MB	2228 MB	89.03	89.00
Amazon2M (3-layer)	1961s	1523s	11218 MB	2235 MB	90.21	90.21
Amazon2M (4-layer)	N/A	2289s	OOM	2241 MB	N/A	90.41

# GraphSAINT

- Unbiased sampler
  - Random node/edge/walk sampler
- Unbiased aggregated representations

$$\zeta_v^{(\ell+1)} = \sum_{u \in \mathcal{V}} \frac{\tilde{A}_{v,u}}{\alpha_{u,v}} \left( \mathbf{W}^{(\ell)} \right)^\top \mathbf{x}_u^{(\ell)} \mathbb{1}_{u|v} = \sum_{u \in \mathcal{V}} \frac{\tilde{A}_{v,u}}{\alpha_{u,v}} \tilde{\mathbf{x}}_u^{(\ell)} \mathbb{1}_{u|v}$$



# GraphSAINT Performance

Dataset	Nodes	Edges	Degree	Feature	Classes	Train / Val / Test
PPI	14,755	225,270	15	50	121 (m)	0.66 / 0.12 / 0.22
Flickr	89,250	899,756	10	500	7 (s)	0.50 / 0.25 / 0.25
Reddit	232,965	11,606,919	50	602	41 (s)	0.66 / 0.10 / 0.24
Yelp	716,847	6,977,410	10	300	100 (m)	0.75 / 0.10 / 0.15
Amazon	1,598,960	132,169,734	83	200	107 (m)	0.85 / 0.05 / 0.10
PPI (large version)	56,944	818,716	14	50	121 (m)	0.79 / 0.11 / 0.10

Method	PPI	Flickr	Reddit	Yelp	Amazon
GCN	0.515±0.006	0.492±0.003	0.933±0.000	0.378±0.001	0.281±0.005
GraphSAGE	0.637±0.006	0.501±0.013	0.953±0.001	0.634±0.006	0.758±0.002
FastGCN	0.513±0.032	0.504±0.001	0.924±0.001	0.265±0.053	0.174±0.021
S-GCN	0.963±0.010	0.482±0.003	0.964±0.001	0.640±0.002	—‡
AS-GCN	0.687±0.012	0.504±0.002	0.958±0.001	—‡	—‡
ClusterGCN	0.875±0.004	0.481±0.005	0.954±0.001	0.609±0.005	0.759±0.008
GraphSAINT-Node	0.960±0.001	0.507±0.001	0.962±0.001	0.641±0.000	0.782±0.004
GraphSAINT-Edge	<b>0.981±0.007</b>	0.510±0.002	<b>0.966±0.001</b>	<b>0.653±0.003</b>	0.807±0.001
GraphSAINT-RW	<b>0.981±0.004</b>	<b>0.511±0.001</b>	<b>0.966±0.001</b>	<b>0.653±0.003</b>	<b>0.815±0.001</b>
GraphSAINT-MRW	0.980±0.006	0.510±0.001	0.964±0.000	0.652±0.001	0.809±0.001

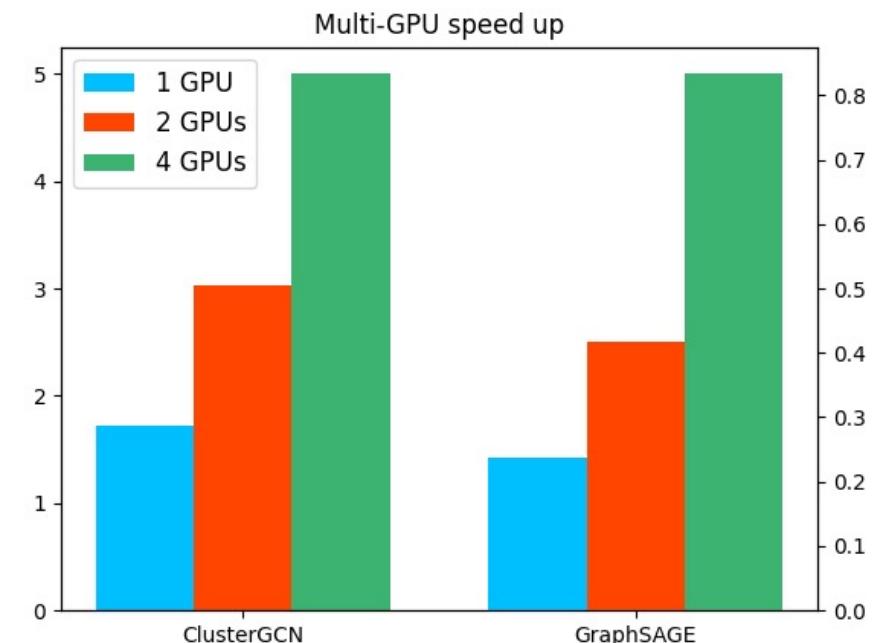
# Multi-GPU Training

Multi-GPU implementation:  
sampling + PyTorch DDP

✓ **Sampler in CogDL**

- [+] NeighborSampling
- [+] ClusterGCN
- [+] GraphSAINT

✓ 4 GPUs ~ 3x↑ speedup



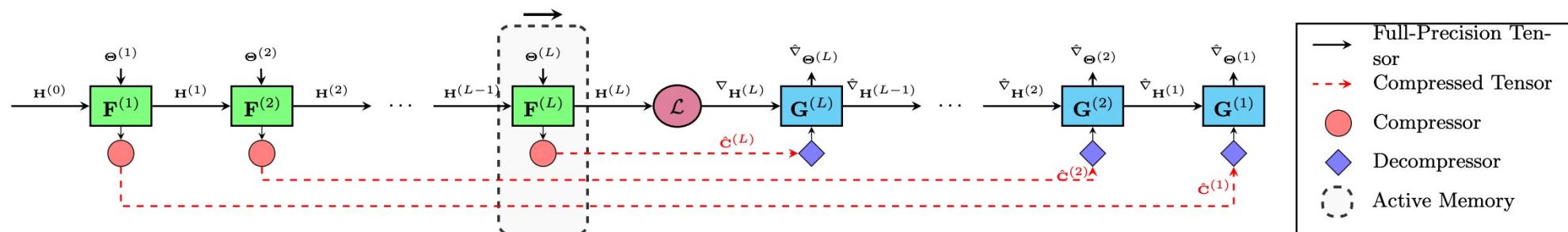
**Usage:** `python scripts/train.py --model gcn --task node_classification  
--dataset reddit --trainer dist_clustergcn`

# Other Solutions for very Deep GNNs?

- GPU memory is the bottleneck for training very deep GNNs.
- Recall RevGNN uses reversible blocks.
- Are there other solutions?
- Activation Compressed Training!

# ActNN : Activation Compressed Training

- ActNN : Reducing Training Memory Footprint via **2-Bit Activation Compressed Training** (By Jianfei Chen, Tsinghua)
- “ActNN reduces the memory footprint of the activation by  $12\times$ .”
- <https://github.com/ucbrise/actnn>



# ActNN Theory

**Theorem 1.** (*Unbiased Gradient*) There exists random quantization strategies for  $\hat{\mathbf{C}}$ , such that

$$\mathbb{E} \left[ \hat{\nabla}_{\Theta} \right] = \nabla_{\Theta} \mathcal{L}_{\mathcal{D}}(\Theta).$$

**Theorem 2.** (*Convergence*) If A1-A3 holds, and  $0 < \alpha \leq \frac{1}{\beta}$ , take the number of iterations  $t$  uniformly from  $\{1, \dots, T\}$ , where  $T$  is a maximum number of iterations. Then

$$\mathbb{E} \|\nabla \mathcal{L}_{\mathcal{D}}(\Theta_t)\|^2 \leq \frac{2(\mathcal{L}(\Theta_1) - \mathcal{L}_{inf})}{\alpha T} + \alpha \beta \sigma^2. \quad (6)$$

**Theorem 3.** (*Gradient Variance*)

$$\text{Var} \left[ \hat{\nabla}_{\Theta^{(l)}} \right] = \text{Var} [\nabla_{\Theta^{(l)}}] + \sum_{m=l}^L \mathbb{E} \left[ \text{Var} \left[ \mathbf{G}_{\Theta}^{(l \sim m)} \left( \hat{\nabla}_{\mathbf{H}^{(m)}}, \hat{\mathbf{C}}^{(m)} \right) \mid \hat{\nabla}_{\mathbf{H}^{(m)}} \right] \right].$$

# ActNN Implementation

```
class RegularLayer:  
    def forward(context, input):  
        context.save_for_backward(input)  
        return compute_output(input)  
  
    def backward(context, grad_output):  
        input = context.saved_tensors  
        return compute_gradient(grad_output, input)  
  
class ActivationCompressedLayer:  
    def forward(context, input):  
        context.save_for_backward(compress(input))  
        return compute_output(input)  
  
    def backward(context, grad_output):  
        input = decompress(context.saved_tensors))  
        return compute_gradient(grad_output, input)
```

# ActNN Performance

- Experiment on ImageNet

Bits	32	4	3	2	1.5	1.25
FP	<b>77.1</b>	N/A	N/A	N/A	N/A	N/A
BLPA	N/A	<b>76.6</b>	Div.	Div.	N/A	N/A
ActNN (L2)	N/A	-	<b>77.4</b>	0.1	N/A	N/A
ActNN (L2.5)	N/A	-	-	<b>77.1</b>	75.9	75.1
ActNN (L3)	N/A	-	-	<b>76.9</b>	76.4	75.9

Network	Batch	Total Mem. (GB)			Act. Mem. (GB)		
		FP	ActNN (L3)	R	FP	ActNN (L3)	R
ResNet-152	32	6.01	1.18	5×	5.28	0.44	12×
	64	11.32	1.64	7×	10.57	0.88	12×
	96	OOM	2.11	/	OOM	1.32	/
	512	OOM	8.27	/	OOM	7.01	/
FCN-HR-48	2	5.76	1.39	4×	4.76	0.39	12×
	4	10.52	1.79	6×	9.52	0.79	12×
	6	OOM	2.17	/	OOM	1.18	/
	20	OOM	4.91	/	OOM	3.91	/

# When SpMM meets ActNN (in CogDL)

$$H^{(i+1)} = AH^{(i)}W$$

```
class SPMMFunction(torch.autograd.Function):
    @staticmethod
    def forward(ctx, rowptr, colind, feat, edge_weight_csr=None, sym=False):
        if edge_weight_csr is None:
            out = spmm.csr_spmm_no_edge_value(rowptr, colind, feat)
        else:
            out = spmm.csr_spmm(rowptr, colind, edge_weight_csr, feat)
        ctx.backward_csc = (rowptr, colind, feat, edge_weight_csr, sym)
        return out

    @staticmethod
    def backward(ctx, grad_out):
        rowptr, colind, feat, edge_weight_csr, sym = ctx.backward_csc
        if edge_weight_csr is not None:
            grad_out = grad_out.contiguous()
            if sym:
                colptr, rowind, edge_weight_csc = rowptr, colind, edge_weight_csr
            else:
                colptr, rowind, edge_weight_csc = spmm.csr2csc(rowptr, colind, edge_weight_csr)
            grad_feat = spmm.csr_spmm(colptr, rowind, edge_weight_csc, grad_out)
            grad_edge_weight = sddmm.csr_sddmm(rowptr, colind, grad_out, feat)
        else:
            if sym is False:
                colptr, rowind, edge_weight_csc = spmm.csr2csc(rowptr, colind, edge_weight_csr)
                grad_feat = spmm.csr_spmm_no_edge_value(colptr, rowind, grad_out)
            else:
                grad_feat = spmm.csr_spmm_no_edge_value(rowptr, colind, grad_out)
                grad_edge_weight = None
        return None, None, grad_feat, grad_edge_weight, None
```

```
class ActSPMMFunction(torch.autograd.Function):
    @staticmethod
    def forward(ctx, rowptr, colind, feat, edge_weight_csr=None, sym=False):
        if edge_weight_csr is None:
            out = spmm.csr_spmm_no_edge_value(rowptr, colind, feat)
        else:
            out = spmm.csr_spmm(rowptr, colind, edge_weight_csr, feat)
        quantized = quantize_activation(feat, None)
        ctx.backward_csc = (rowptr, colind, quantized, edge_weight_csr, sym)
        ctx.other_args = feat.shape
        return out

    @staticmethod
    def backward(ctx, grad_out):
        rowptr, colind, quantized, edge_weight_csr, sym = ctx.backward_csc
        q_input_shape = ctx.other_args
        feat = dequantize_activation(quantized, q_input_shape)
        del quantized, ctx.backward_csc

        if edge_weight_csr is not None:
            grad_out = grad_out.contiguous()
            if sym:
                colptr, rowind, edge_weight_csc = rowptr, colind, edge_weight_csr
            else:
                colptr, rowind, edge_weight_csc = spmm.csr2csc(rowptr, colind, edge_weight_csr)
            grad_feat = spmm.csr_spmm(colptr, rowind, edge_weight_csc, grad_out)
            grad_edge_weight = sddmm.csr_sddmm(rowptr, colind, grad_out, feat)
```

# Experimental Results

- Default setting of CogDL

Dataset	Origin GCN	GCN + actnn
Cora	$81.30 \pm 0.22$	$81.27 \pm 0.19$
Citeseer	$71.73 \pm 0.54$	$71.70 \pm 0.28$
Pubmed	$79.17 \pm 0.12$	$79.10 \pm 0.08$
Flickr	$50.74 \pm 0.10$	$50.89 \pm 0.04$
Reddit	$95.01 \pm 0.02$	$94.89 \pm 0.01$

# Activation Memory (GCN + ActNN)

- Setting:  $H = \text{Dropout}(\text{ReLU}(\text{BN}(AHW)))$

#dataset, #layers, #hidden	Origin GCN	GCN + actnn	ratio	Ideal ratio
PPI, 5, 2048	3704	420	8.8x	raw: $32*2 / (2.125*2+1) = 12.2x$ +bn: $32*3 / (2.125*3+1) = 13.0x$ +bn+dropout: $32*4 / (2.125*3+2) = 15.3x$
PPI, 5, 2048 (+bn)	5484	539	10.2x	
PPI, 5, 2048 (+bn, +dropout)	7711	594	13.0x	
Flickr, 5, 512	1420	154	9.2x	
Flickr, 5, 512 (+bn)	2117	201	10.5x	
Flickr, 5, 512 (+bn, +dropout)	2991	223	13.4x	
Flickr, 10, 512	3178	311	10.2x	
Flickr, 10, 512 (+bn)	4747	415	11.4x	
Flickr, 10, 512 (+bn, +dropout)	6712	465	14.4x	

# Activation Memory (GraphSAGE + ActNN)

- Setting:  $H = \text{Dropout} \left( \text{ReLU} \left( \text{BN} \left( \text{Concat}(AH, H)W \right) \right) \right)$

#dataset, #layers, #hidden	Origin SAGE	SAGE + actnn	ratio	Ideal ratio
PPI, 5, 2048	5524	580	9.5x	raw: $32*2 / (2.125*2+1) = 12.2x$ +bn: $32*3 / (2.125*3+1) = 13.0x$ +bn+dropout: $32*4 / (2.125*3+2) = 15.3x$
PPI, 5, 2048 (+bn)	7304	698	10.5x	
PPI, 5, 2048 (+bn, +dropout)	OOM	754	-	
Flickr, 5, 512	2457	209	11.8x	
Flickr, 5, 512 (+bn)	3155	255	12.4x	
Flickr, 5, 512 (+bn, +dropout)	4027	278	14.5x	
Flickr, 10, 512	5090	430	11.8x	
Flickr, 10, 512 (+bn)	6659	534	12.5x	
Flickr, 10, 512 (+bn, +dropout)	8624	584	14.8x	

# Activation Memory (GCNII + ActNN)

- Setting:

$$\mathbf{H}^{(\ell+1)} = \sigma \left( \left( (1 - \alpha_\ell) \tilde{\mathbf{P}} \mathbf{H}^{(\ell)} + \alpha_\ell \mathbf{H}^{(0)} \right) \left( (1 - \beta_\ell) \mathbf{I}_n + \beta_\ell \mathbf{W}^{(\ell)} \right) \right)$$

#dataset, #layers, #hidden	GCNII	GCNII + actnn	ratio	ideal
PPI, 10, 512	4008	340	11.8x	$(32*3) / (2*2.125+2) = 15.36$
PPI, 20, 512	7708	619	12.5x	
PPI, 20, 1024	7879	603	13.1x	
Flickr, 5, 512	3421	229	14.9x	
Flickr, 10, 512	6268	413	15.2x	
Flickr, 5, 1024	6660	438	15.2x	

#layer s	Origin GCNII	GCNII + actnn
32	$84.83 \pm 0.33$	$84.67 \pm 0.12$
64	$85.13 \pm 0.61$	$85.00 \pm 0.37$
128	$84.83 \pm 0.58$	$85.20 \pm 0.36$
256	$85.00 \pm 0.08$	$85.37 \pm 0.54$

# Activation Memory (GIN + ActNN)

- Setting:  $\mathbf{H}^{(l)} = \text{MLP}^{(l)} \left( (1 + \epsilon) \mathbf{H}^{(l-1)} + \mathbf{A} \mathbf{H}^{(l-1)} \right)$   
 $\mathbf{h}_G = \text{CONCAT}(\text{READOUT}(\mathbf{H}^{(l)}), l = 0, 1, \dots L$

#dataset, #batch, #layers, #hidden	GIN	GIN + actnn	ratio	ideal
NCI1, 512, 20, 512	2723	262	10.4x	$(32*3) / (2.125*3+1) =$ 13.0x
NCI1, 512, 20, 1024	5735	540	10.6x	
NCI1, 1024, 20, 512	5502	528	10.4x	
NCI1, 512, 40, 512	6231	590	10.6x	

# Self-supervised Learning on Graphs

- Types of self-supervision:
  - generative learning
  - contrastive learning
- Learning paradigm:
  - Pre-training & Fine-tuning
  - Joint learning
  - Self-training
- Encoders: GCN, GAT, GIN
- Downstream tasks

# Survey of Graph Self-supervised Learning

Table 1: An overview of recent generative models on graphs. For acronyms used, "PT+FT" refers to Pre-training and Fine-tuning, "JL" refers to Joint Learning, "ST" refers to Self-Training.

Model	Self-Supervision Tasks	Training Scheme	Graph Encoders	Downstream Tasks
GAE/VGAE[1]	Graph Reconstruction	PT+FT	GCN	Node/Link
EdgeMask[2]	Graph Reconstruction	JL/PT+FT	GCN	Node
AttributeMask[2]	Graph Reconstruction	JL/PT+FT	GCN	Node
GraphCompletion[3]	Graph Reconstruction	JL/PT+FT	GCN/GAT/GIN/GraphMix	Node
GPT-GNN[4]	Graph Reconstruction	PT+FT	Heterogeneous Graph Transformer	Node/Link
S <sup>2</sup> GRL[5]	Graph Property Prediction	PT+FT	GCN/GraphSAGE	Node
Distance2Cluster[2]	Graph Property Prediction	JL/PT+FT	GCN	Node
Node clustering[3]	Graph Property Prediction	JL/PT+FT	GCN/GAT/GIN/GraphMix	Node
Graph partitioning[3]	Graph Property Prediction	JL/PT+FT	GCN/GAT/GIN/GraphMix	Node
SuperGAT[6]	Graph Property Prediction	JL	SuperGAT	Node
Un_GraphSAGE[7]	Graph Property Prediction	PT+FT	GraphSAGE	Node
M3S[8]	Graph Property Prediction	ST	GCN	Node
AdvT[3]	Adversarial Attack and Defense	JL	GCN	Node
Graph-Bert[9]	Hybrid	PT+FT	Graph Transformer	Node

Table 2: An overview of recent contrastive models on graphs.

Model	Augmentation Method	Contrastive Framework	Graph Encoders	Downstream Tasks
InfoGraph[10]	Original Graph	Deep Infomax	GIN	Graph
GraphCL[11]	Graph Corruption		GCN	Graph
GRACE[12]	Graph Corruption	InfoNCE	GCN	Node
GCC[13]	Random walk Sampling	MoCo	GIN	Node/Graph
DGI[14]	Graph Corruption + Sampling	Deep Infomax	GCN	Node
MVGRL[15]	Graph Corruption + Sampling	Deep Infomax	GCN/GraphSAGE	Node
DwGCL[16]	Graph Corruption + Sampling		GCN	Node

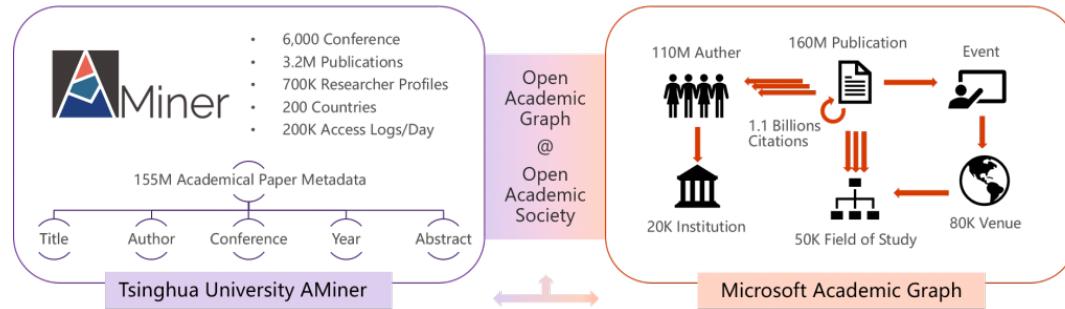
# Results of Self-supervised Learning

- Learning paradigm:
  - Self-supervised (SL), Joint Learning (JL), unsupervised representation learning (URL)
- Semi-supervised datasets: Cora, Citeseer, PubMed
- Supervised datasets: Flickr, Reddit

Model	Training Scheme	Cora	Citeseer	PubMed	Flickr	Reddit
Linear	SL	47.86±0.02	49.25±0.04	69.17±0.03	45.81±0.00	67.95±0.01
GCN	SL	81.53± 0.26	71.75±0.05	79.30±0.31	52.74±0.13	95.16±0.01
EdgeMask	JL	81.28±0.31	71.53±0.24	79.55±0.11	52.55±0.06	95.07±0.01
AttributeMask	JL	81.20±0.32	71.45±0.35	78.93±0.25	52.45±0.15	95.00±0.01
S <sup>2</sup> GRL	JL	83.42±0.63	72.37±0.11	81.20±0.37	52.59±0.05	95.17±0.00
Distance2Clusters	JL	82.47±0.48	71.55±0.30	81.53±0.22	52.24±0.07	/
SuperGAT	JL	82.77±0.53	72.25±0.52	80.30±0.31	52.80±0.07	95.42±0.01
MVGRL	URL	80.76±0.80	65.84±2.72	76.01±2.46	46.08±0.39	92.76±0.22
DGI	URL	81.91±0.17	70.01±0.87	76.49±1.04	46.35±0.11	93.12±0.18
EdgeMask	URL	75.23±1.14	68.96±0.96	79.41±1.15	50.48±0.06	93.68±0.01
AttributeMask	URL	76.12±0.91	70.69±0.44	75.16±1.71	51.57±0.15	93.37±0.01
S <sup>2</sup> GRL	URL	81.56±0.49	69.48±0.91	80.83±0.57	50.85±0.03	93.88±0.05
Distance2Cluster	URL	73.86±0.29	66.53±0.29	79.44±0.34	50.24±0.07	/

# OAG: Open Academic Graph

<https://www.openacademic.ai/oag/>



Data set	#Pairs/Venues	Date
Linking relations	29,841	2018.12
AMiner venues	69,397	2018.07
MAG venues	52,678	2018.11

Table 1: statistics of OAG venue data

Data set	#Pairs/Papers	Date
Linking relations	91,137,597	2018.12
AMiner papers	172,209,563	2019.01
MAG papers	208,915,369	2018.11

Table 2: statistics of OAG paper data

Data set	#Pairs/Authors	Date
Linking relations	1,717,680	2019.01
AMiner authors	113,171,945	2018.07
MAG authors	253,144,301	2018.11

## Open Academic Graph

Open Academic Graph (OAG) is a large knowledge graph unifying two billion-scale academic graphs: Microsoft Academic Graph (MAG) and AMiner. In mid 2017, we published OAG v1, which contains 166,192,182 papers from MAG and 154,771,162 papers from AMiner (see below) and generated 64,639,608 linking (matching) relations between the two graphs. This time, in OAG v2, author, venue and newer publication data and the corresponding matchings are available.

## Overview of OAG v2

The statistics of OAG v2 is listed as the three tables below. The two large graphs are both evolving and we take MAG November 2018 snapshot and AMiner July 2018 or January 2019 snapshot for this version.

# Heterogeneous Graph Benchmark (HGB)

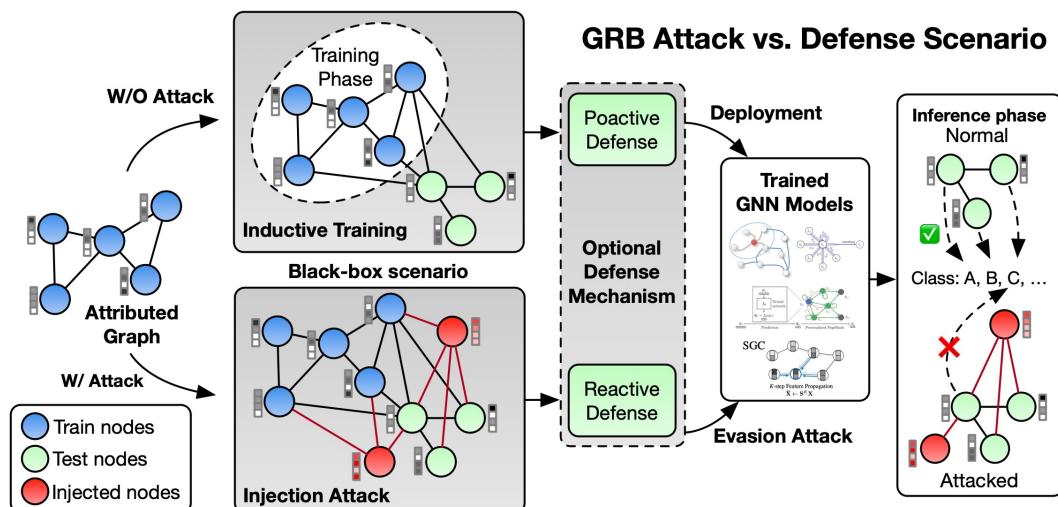
- A unified benchmark datasets and evaluation pipelines for heterogeneous graph research.
- **Paper:** Are we really making much progress? Revisiting, benchmarking and refining heterogeneous graph neural networks. (*KDD'21*)
- **Code & Data:** <https://github.com/THUDM/HGB>
- **Leaderboard:** <https://www.biendata.xyz/hgb/>
- There is also a simple baseline Simple-HGN in HGB. We find that a rather simple design of heterogeneous GNN can reach SOTA.

## Background:

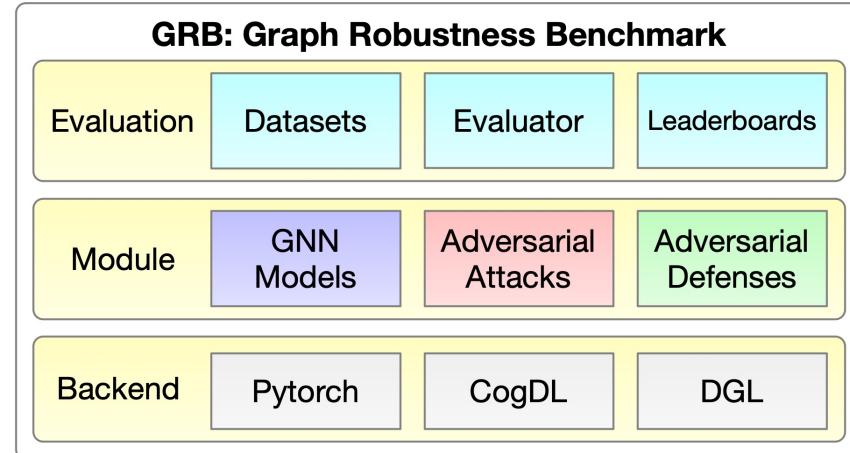
Recently, works have proved that adversarial attacks can threaten the *robustness* of graph ML models in various tasks.

## Problems:

1. Ill-defined threat model in previous works.
2. Absence of unified and standard evaluation approach.



Example of GRB evaluation scenario



GRB framework

## Solution: Graph Robustness Benchmark (GRB)

*Scalable, general, unified, and reproducible* benchmark on *adversarial robustness* of graph ML models, which facilitates fair comparisons among various attacks & defenses and promotes future research in this field.

*Graph Robustness Benchmark: Rethinking and Benchmarking Adversarial Robustness of Graph Neural Networks*

Qinkai Zheng, Xu Zou, Yuxiao Dong, Yukuo Cen, Jie Tang



*All discussions and contributions  
are highly welcome!*

Graph Robustness Benchmark (GRB): Key Features			
<b>Elaborated Datasets</b>	<b>Scalability</b>	Datasets from small to large scales.	
	<b>Specificity</b>	Novel splitting scheme + Preprocessing.	
<b>Modular Framework</b>	<b>Implementation</b>	GNNs (GCN, SAGE, GAT, GIN, APPNP, ...)	
		Attacks (FGSM, PGD, SPEIT, TDGIA, ...)	
		Defenses (Adversarial Training, GNNGuard, ...)	
<b>Unified Evaluation</b>	<b>Backend</b>	Pytorch	CogDL
	<b>Scenario</b>	Well-defined realistic threat model.	
	<b>Fairness</b>	Unified settings for attackers and defenders.	
<b>Reproducible Leaderboard</b>	<b>Pipeline</b>	Easy-to-use evaluation pipeline.	
	<b>Reproducibility</b>	Availability of methods and related materials.	
	<b>Up-to-date</b>	Continuously maintained to track progress.	

**Homepage:**

<https://cogdl.ai/grb/home>

**Github:**

<https://github.com/THUDM/grb>

**Leaderboard:**

<https://cogdl.ai/grb/leaderboard/>

**Docs:** <https://grb.readthedocs.io/>

**Google Group:**

<https://groups.google.com/g/graph-robustness-benchmark>

**Contact:**

[cogdl.grbteam@gmail.com](mailto:cogdl.grbteam@gmail.com)

[qinkai.zheng1028@gmail.com](mailto:qinkai.zheng1028@gmail.com)

# Open Graph Benchmark

- Large-scale, realistic, and diverse benchmark datasets for graph ML.



Paper: <https://arxiv.org/abs/2005.00687>

Webpage: <https://ogb.stanford.edu/>

Github: <https://github.com/snap-stanford/ogb>

# Recommendation Application

- Build recommendation via pipeline API
- Integrate LightGCN (SIGIR'20)
- Similar to Amazon Personalize

```
import numpy as np
from cogdl import pipeline

data = np.array([[0, 0], [0, 1], [0, 2], [1, 1], [1, 3], [1, 4], [2, 4], [2, 5], [2, 6]])
rec = pipeline("recommendation", model="lightgcn", data=data, max_epoch=1)
print(rec([0]))

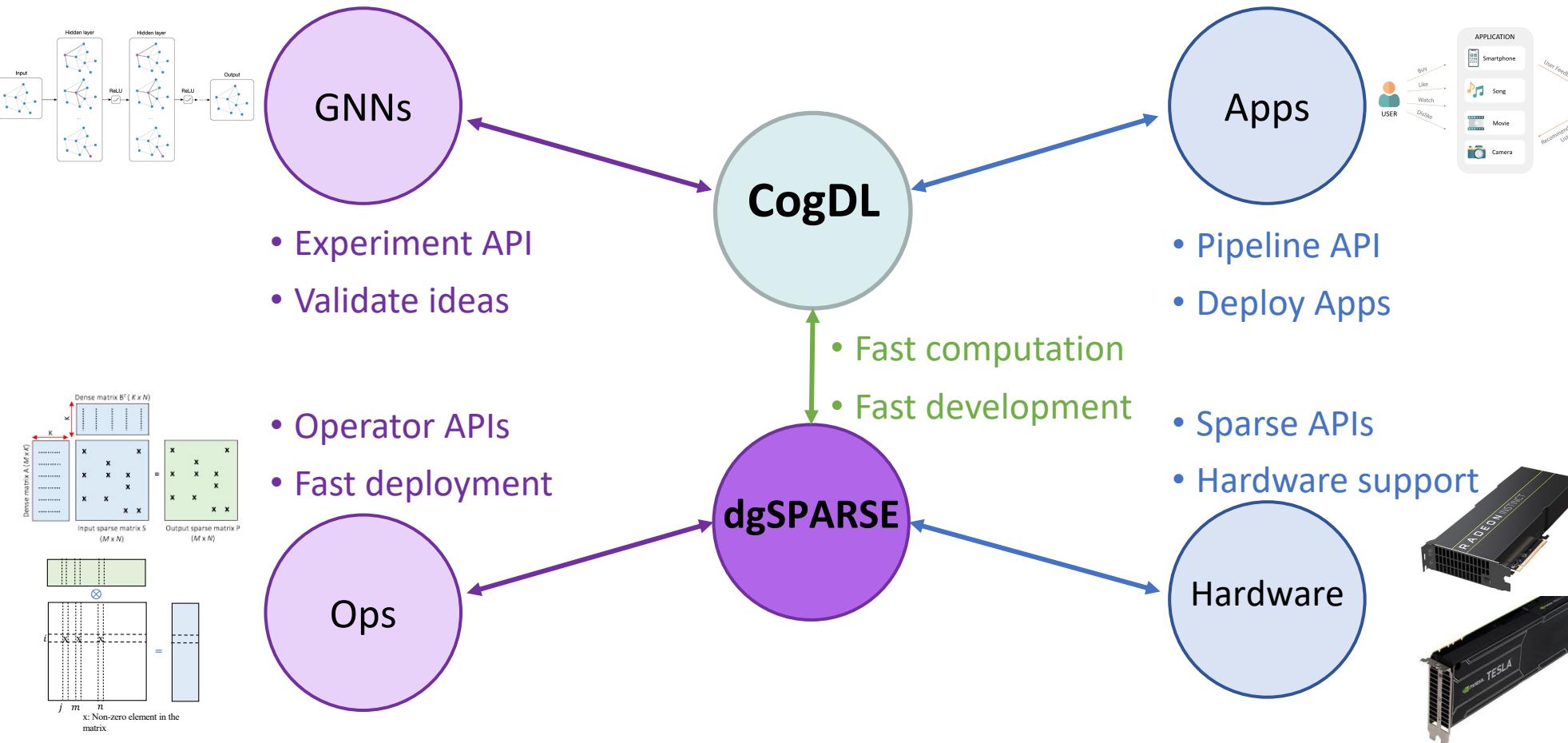
rec = pipeline("recommendation", model="lightgcn", dataset="ali", max_epoch=1)
print(rec([0]))
```

# AMiner Subscribe

- Recommend papers, scholars to users

The screenshot shows the AMiner Research Feed page. On the left, there's a sidebar with 'Profile Management' (Academic Profile, User Information), 'Academic Home' (Research Feed, My Following, Paper Collections), and a 'Research Feed' section which is currently selected. The main content area has a header 'What can I do here?' with four options: 'Track the Latest Paper' (checked), 'Claim My Academic Homepage' (unchecked), 'Promote My Paper' (unchecked), and 'Receive the latest news' (checked). Below this, it says 'Added interest area: Reinforcement Learning, Deep Learning, Knowledge Graph, social network system'. It then displays two recommended scholars: David Silver and Pieter Abbeel, each with a profile picture, name, h-index, #Paper, #Citation, and a 'Follow' button. Below them is a section titled 'High Quality Paper List on Conference or Topic' for 'CVPR2020', showing a thumbnail, title, authors, and a link to the paper. At the bottom, there are two more paper entries: 'Finite-Sample Analysis for SARSA with Linear Function Approximation' by Shaofeng Zou, Tengyu Xu, Yingbin Liang, and 'Tighter Problem-Dependent Regret Bounds in Reinforcement Learning without Domain Knowledge using Value Function Bounds' by Andrea Zanette, Emma Brunskill.

# CogDL & dgSPARSE



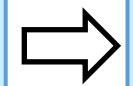
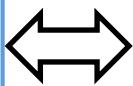
# GRL: NE&GNN

Network  
Embedding

Matrix  
Factorization

Graph Neural  
Networks

GNN  
Pre-Training



Learning GNNs with **CogDL**



<https://alchemy.tencent.com/>



<https://github.com/thudm/cogdl>

# Related Publications

For more, check <http://keg.cs.tsinghua.edu.cn/jietang>

- Jialin Zhao, Yuxiao Dong, Ming Ding, Evgeny Kharlamov, and Jie Tang. Adaptive Diffusion in Graph Neural Networks. **NeurIPS'21**.
- Wenzheng Feng, Jie Zhang, Yuxiao Dong, Yu Han, Huanbo Luan, Qian Xu, Qiang Yang, Evgeny Kharlamov, and Jie Tang. Graph Random Neural Networks for Semi-Supervised Learning on Graphs. **NeurIPS'20**.
- Ming Ding, Chang Zhou, Hongxia Yang, and Jie Tang. CogLTX: Applying BERT to Long Texts. **NeurIPS'20**.
- Jiezhong Qiu, Chi Wang, Ben Liao, Richard Peng, and Jie Tang. Concentration Bounds for Co-occurrence Matrices of Markov Chains. **NeurIPS'20**.
- Xu Zou, Qinkai Zheng, Yuxiao Dong, Xinyu Guan, Evgeny Kharlamov, Jialiang Lu, and Jie Tang. TDGIA: Effective Injection Attacks on Graph Neural Networks. **KDD'21**.
- Qingsong Lv, Ming Ding, Qiang Liu, Yuxiang Chen, Wenzheng Feng, Siming He, Chang Zhou, Jian-guo Jiang, Yuxiao Dong, and Jie Tang. Are we really making much progress? Revisiting, benchmarking and refining the Heterogeneous Graph Neural Networks. **KDD'21**.
- Tinglin Huang, Yuxiao Dong, Ming Ding, Zhen Yang, Wenzheng Feng, Xinyu Wang, and Jie Tang. MixGCF: An Improved Training Method for Graph Neural Network-based Recommender Systems. **KDD'21**.
- Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. GCC: Graph Contrastive Coding for Structural Graph Representation Pre-Training. **KDD'20**.
- Zhen Yang, Ming Ding, Chang Zhou, Hongxia Yang, Jingren Zhou, and Jie Tang. Understanding Negative Sampling in Graph Representation Learning. **KDD'20**.
- Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. Controllable Multi-Interest Framework for Recommendation. **KDD'20**.
- Yukuo Cen, Xu Zou, Jianwei Zhang, Hongxia Yang, Jingren Zhou and Jie Tang. Representation Learning for Attributed Multiplex Heterogeneous Network. **KDD'19**.
- Fanjin Zhang, Xiao Liu, Jie Tang, Yuxiao Dong, Peiran Yao, Jie Zhang, Xiaotao Gu, Yan Wang, Bin Shao, Rui Li, and Kuansan Wang. OAG: Toward Linking Large-scale Heterogeneous Entity Graphs. **KDD'19**.
- Qibin Chen, Junyang Lin, Yichang Zhang, Hongxia Yang, Jingren Zhou and Jie Tang. Towards Knowledge-Based Personalized Product Description Generation in E-commerce. **KDD'19**.
- Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. DeepInf: Modeling Influence Locality in Large Social Networks. **KDD'18**.
- Jiezhong Qiu, Laxman Dhulipala, Jie Tang, Richard Peng, and Chi Wang. LightNE: A Lightweight Graph Processing System for Network Embedding. **SIGMOD'21**.
- Yuxiao Dong, Ziniu Hu, Kuansan Wang, Yizhou Sun and Jie Tang. Heterogeneous Network Representation Learning. **IJCAI'20**.
- Jie Zhang, Yuxiao Dong, Yan Wang, Jie Tang, and Ming Ding. ProNE: Fast and Scalable Network Representation Learning. **IJCAI'19**.
- Yifeng Zhao, Xiangwei Wang, Hongxia Yang, Le Song, and Jie Tang. Large Scale Evolving Graphs with Burst Detection. **IJCAI'19**.
- Yu Han, Jie Tang, and Qian Chen. Network Embedding under Partial Monitoring for Evolving Networks. **IJCAI'19**.
- Yifeng Zhao, Xiangwei Wang, Hongxia Yang, Le Song, and Jie Tang. Large Scale Evolving Graphs with Burst Detection. **IJCAI'19**.
- Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. Cognitive Graph for Multi-Hop Reading Comprehension at Scale. **ACL'19**.
- Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Chi Wang, Kuansan Wang, and Jie Tang. NetSMF: Large-Scale Network Embedding as Sparse Matrix Factorization. **WWW'19**.
- Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec. **WSDM'18**.
- Zhenyu Hou, Yukuo Cen, Yuxiao Dong, Jie Zhang, and Jie Tang. Automated Unsupervised Graph Representation Learning. **TKDE'21**.
- Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-supervised Learning: Generative or Contrastive. **TKDE'21**.
- Shu Zhao, Ziwei Du, Jie Chen, Yanping Zhang, Jie Tang, and Philip S. Yu. Hierarchical Representation Learning for Attributed Networks. **TKDE'21**.

# Thank you !

## Collaborators:

Jie Zhang, Ming Ding, Jiezhong Qiu, Qibin Chen, Yifeng Zhao, Yukuo Cen, Yu Han, Fanjin Zhang, Xu Zou, Yan Wang, et al. (**THU**)

Yuxiao Dong, Kuansan Wang (**Microsoft**)

Hongxiao Yang, Chang Zhou, Le Song, Jingren Zhou, et al. (**Alibaba**)

Jie Tang, KEG, Tsinghua U  
Download all data & Codes

<http://keg.cs.tsinghua.edu.cn/jietang>  
<https://keg.cs.tsinghua.edu.cn/cogdl/>  
<https://github.com/THUDM>